

Unsupervised Learning

Instructor: Alan Ritter

Slides Adapted from Carlos Guestrin Dan Klein and Luke Zettlemoyer

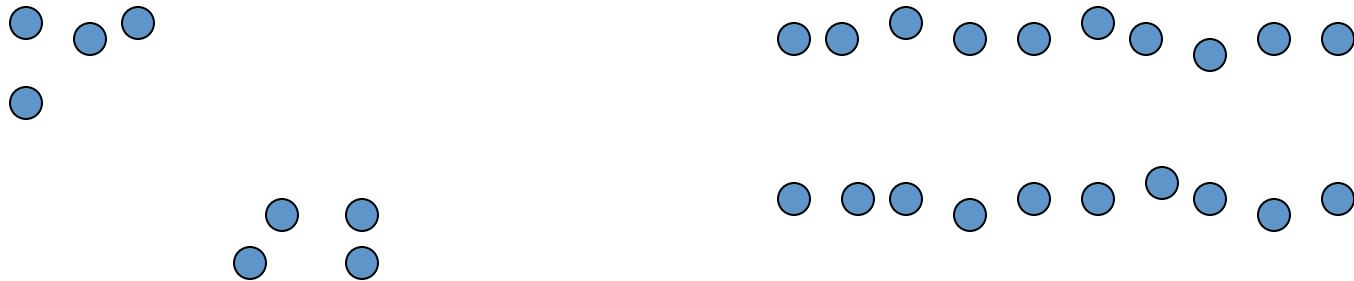
Clustering

- Clustering systems:
 - Unsupervised learning
 - Detect patterns in unlabeled data
 - E.g. group emails or search results
 - E.g. find categories of customers
 - E.g. detect anomalous program executions
 - Useful when don't know what you're looking for
 - Requires data, but no labels
 - Often get gibberish



Clustering

- Basic idea: group together similar instances
- Example: 2D point patterns

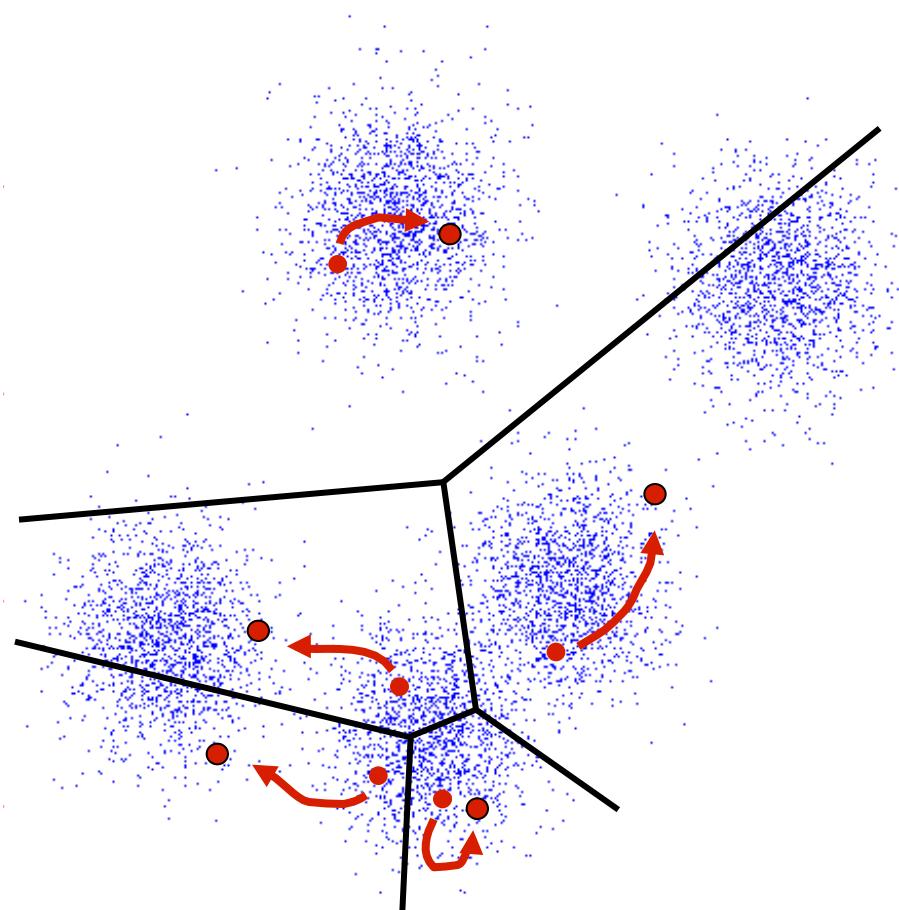


- What could “similar” mean?
 - One option: small (squared) Euclidean distance

$$dist(x, x') = \|x - x'\|_2^2 = (x - x')^T (x - x') = \sum_i (x_i - x'_i)^2$$

K-Means

- An iterative clustering algorithm
 - Pick K random points as cluster centers (means), $c^1 \dots c^k$
 - Alternate:
 - Assign each example x^i to the mean c^j that is closest to it
 - Set each mean c^j to the average of its assigned points
 - Stop when no points' assignments change



Example: K-Means for Segmentation

$K = 2$



$K = 3$



$K = 10$



Original image



K-Means

- Data: $\{x^j \mid j=1..n\}$
- An iterative clustering algorithm
 - Pick K random cluster centers, $c^1 \dots c^k$
 - For $t=1..T$: [or, stop if assignments don't change]
 - for $j = 1..n$: [recompute cluster assignments]

$$a^j = \arg \min_i dist(x^j, c^i)$$

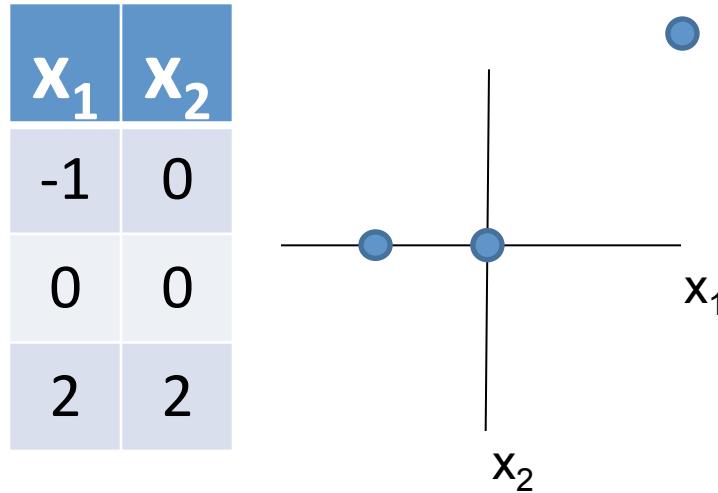
- for $j=1..k$: [recompute cluster centers]

$$c^j = \frac{1}{|\{i | a^i = j\}|} \sum_{\{i | a^i = j\}} x^i$$

Pick K random cluster centers, $c^1 \dots c^k$

For $t=1..T$:

- for $j = 1..n$: [recompute assignments]
 $a^j = \arg \min_i dist(x^j, c^i)$
- for $j = 1..k$: [recompute cluster centers]
 $c^j = \frac{1}{|\{i|a^i=j\}|} \sum_{\{i|a^i=j\}} x^i$



$$dist(x, x') = \sum_i (x_i - x'_i)^2$$

Random cluster means:

- $c^1 = [-1, 0]$, $c^2 = [0, 0]$

$t=0$:

$d(x^j, c^i)$	x^1	x^2	x^3
c^1	0	1	13
c^2	1	0	8

- $a^1 = \operatorname{argmin}_i dist(x^1, c^i) = 1$
- $a^2 = \operatorname{argmin}_i dist(x^2, c^i) = 2$
- $a^3 = \operatorname{argmin}_i dist(x^3, c^i) = 2$
- $c^1 = (1/1) * [-1, 0] = [-1, 0]$
- $c^2 = (1/2) * ([0, 0] + [2, 2]) = [1, 1]$

$t=1$:

$d(x^j, c^i)$	x^1	x^2	x^3
c^1	0	1	13
c^2	4	4	18

- $a^1 = \operatorname{argmin}_i dist(x^1, c^i) = 1$
- $a^2 = \operatorname{argmin}_i dist(x^2, c^i) = 1$
- $a^3 = \operatorname{argmin}_i dist(x^3, c^i) = 2$
- $c^1 = (1/2) * ([-1, 0] + [0, 0]) = [-0.5, 0]$
- $c^2 = (1/1) * [2, 2] = [2, 1]$

$t=2$:

- Stop!! (cluster assignments a^i won't change in next round; you can verify!)

K-Means as Optimization

- Consider the total distance to the means:

$$L(\{x^i\}, \{a^j\}, \{c^k\}) = \sum_i \text{dist}(x^i, c^{a^i})$$

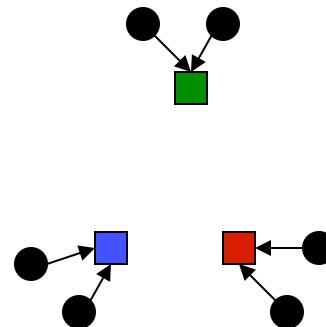
points assignments means

- Two stages each iteration:
 - Update assignments: fix means c , change assignments a
 - Update means: fix assignments a , change means c

- Coordinate gradient descent on L*

- Will it converge?

- Yes!, if you can argue that each update can't increase L



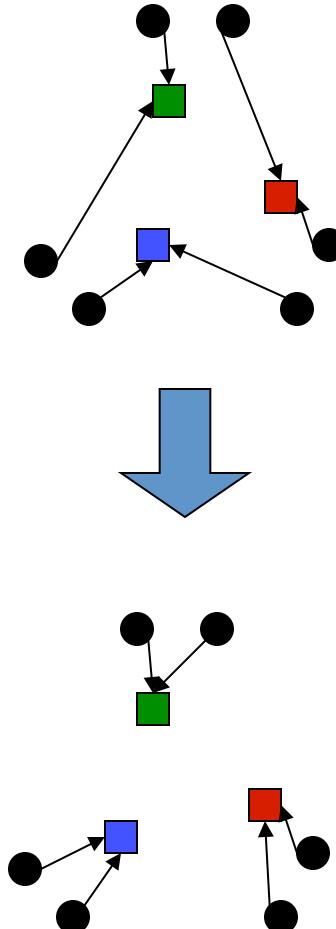
Phase I: Update Assignments

- For each point, re-assign to closest mean:

$$a^j = \arg \min_i dist(x^j, c^i)$$

- Can only decrease total distance L!

$$L(\{x^i\}, \{a^j\}, \{c^k\}) = \sum_i dist(x^i, c^{a^i})$$

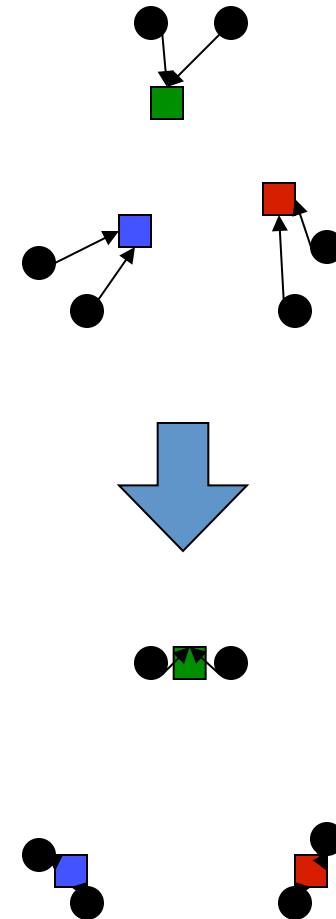


Phase II: Update Means

- Move each mean to the average of its assigned points:

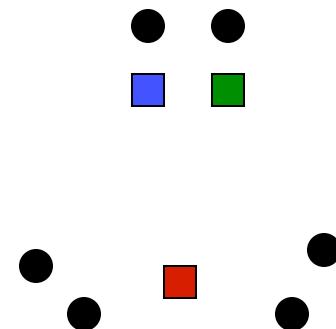
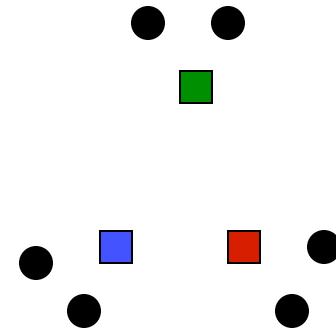
$$c^j = \frac{1}{|\{i|a^i = j\}|} \sum_{\{i|a^i=j\}} x^i$$

- Also can only decrease total distance... (Why?)
- Fun fact: the point y with minimum squared Euclidean distance to a set of points $\{x\}$ is their mean



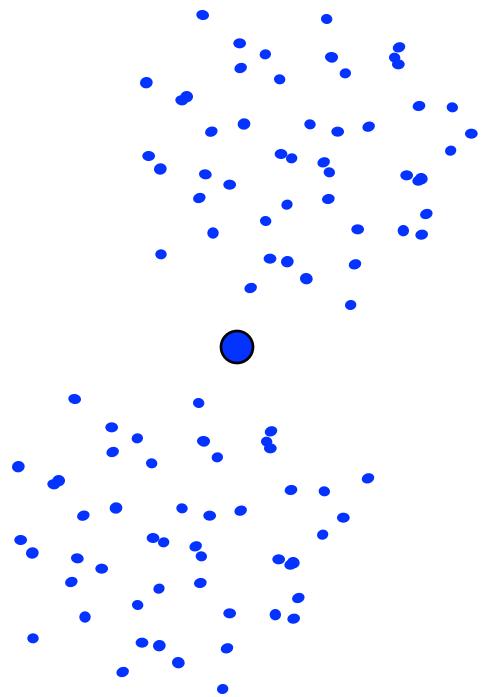
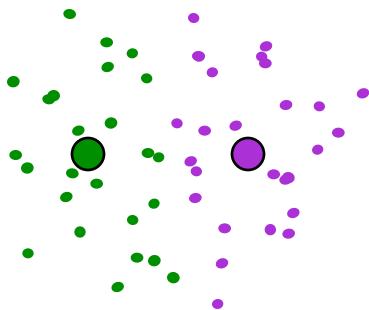
Initialization

- K-means is non-deterministic
 - Requires initial means
 - It does matter what you pick!
 - What can go wrong?
 - Various schemes for preventing this kind of thing: variance-based split / merge, initialization heuristics



K-Means Getting Stuck

- A local optimum:



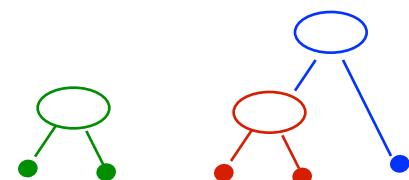
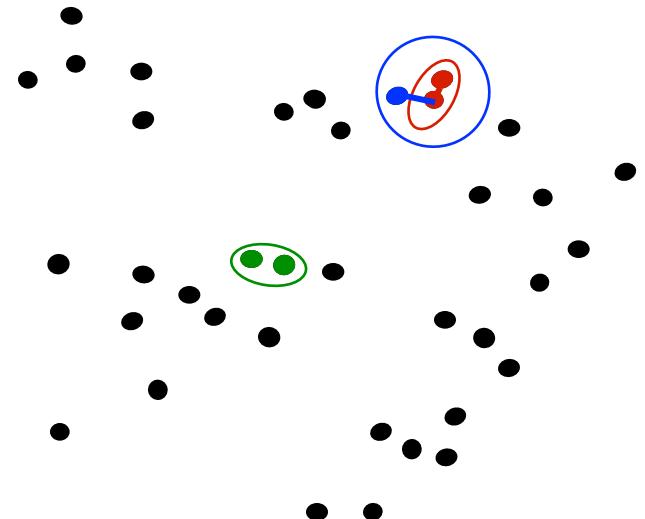
*Why doesn't this work out like
the earlier example, with the
purple taking over half the blue?*

K-Means Questions

- Will K-means converge?
 - To a global optimum?
- Will it always find the true patterns in the data?
 - If the patterns are very very clear?
- Will it find something interesting?
- Do people ever use it?
- How many clusters to pick?

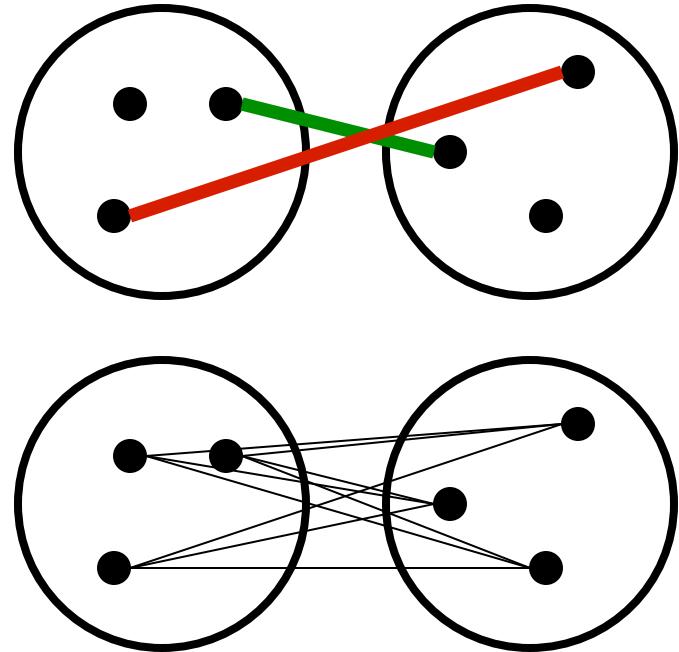
Agglomerative Clustering

- Agglomerative clustering:
 - First merge very similar instances
 - Incrementally build larger clusters out of smaller clusters
- Algorithm:
 - Maintain a set of clusters
 - Initially, each instance in its own cluster
 - Repeat:
 - Pick the two **closest** clusters
 - Merge them into a new cluster
 - Stop when there's only one cluster left
- Produces not one clustering, but a family of clusterings represented by a **dendrogram**

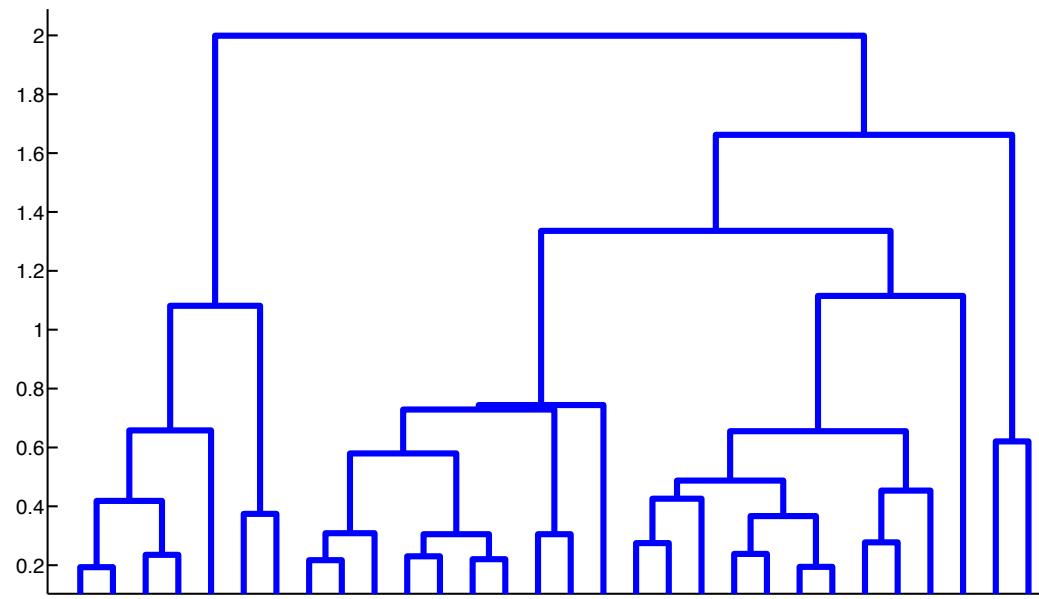


Agglomerative Clustering

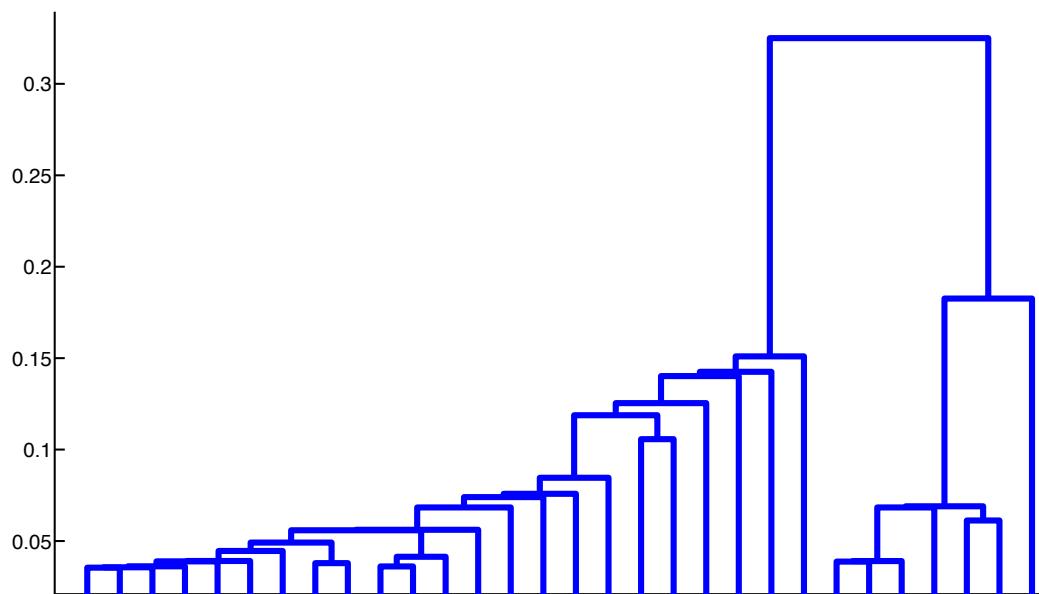
- How should we define “closest” for clusters with multiple elements?
- Many options:
 - Closest pair (single-link clustering)
 - Farthest pair (complete-link clustering)
 - Average of all pairs
- Different choices create different clustering behaviors



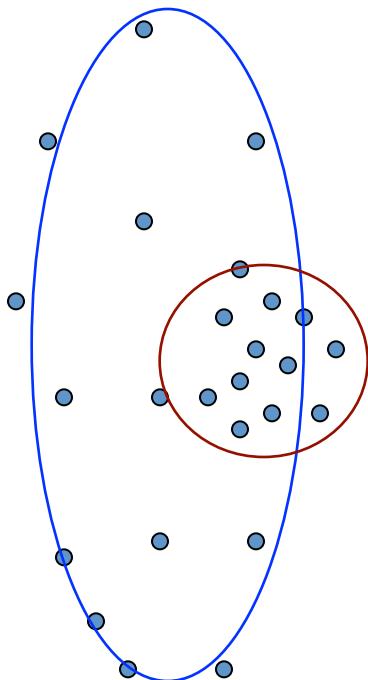
complete link



single link

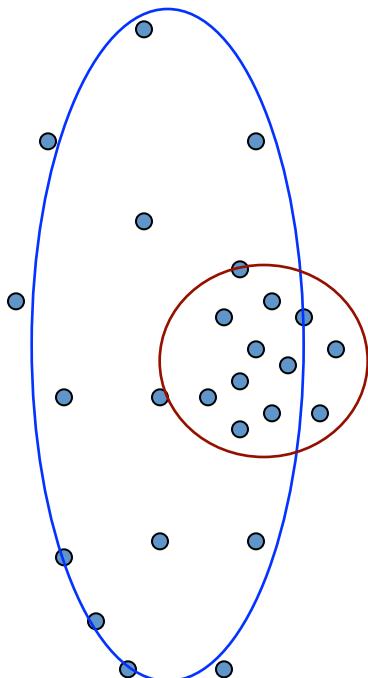


(One) bad case for “hard assignments”?



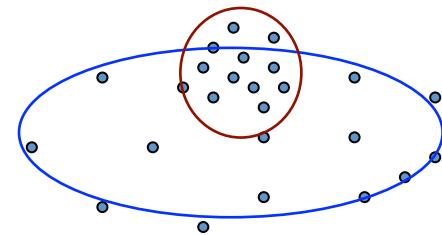
- Clusters may overlap
- Some clusters may be “wider” than others

(One) bad case for “hard assignments”?



- Clusters may overlap
- Some clusters may be “wider” than others
- Distances can be deceiving!

Probabilistic Clustering



- We can use a probabilistic model!
 - allows overlaps, clusters of different size, etc.
- Can tell a *generative story* for data
 - $P(X|Y)$ $P(Y)$ is common
- Challenge: we need to estimate model parameters without labeled Ys

Y	X ₁	X ₂
??	0.1	2.1
??	0.5	-1.1
??	0.0	3.0
??	-0.1	-2.0
??	0.2	1.5
...

What Model Should We Use?

- Depends on X!
- Here, maybe Gaussian Naïve Bayes?
 - Multinomial over clusters Y, Gaussian over each X_i given Y

$$p(Y_i = y_k) = \theta_k$$

$$P(X_i = x \mid Y = y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{\frac{-(x-\mu_{ik})^2}{2\sigma_{ik}^2}}$$

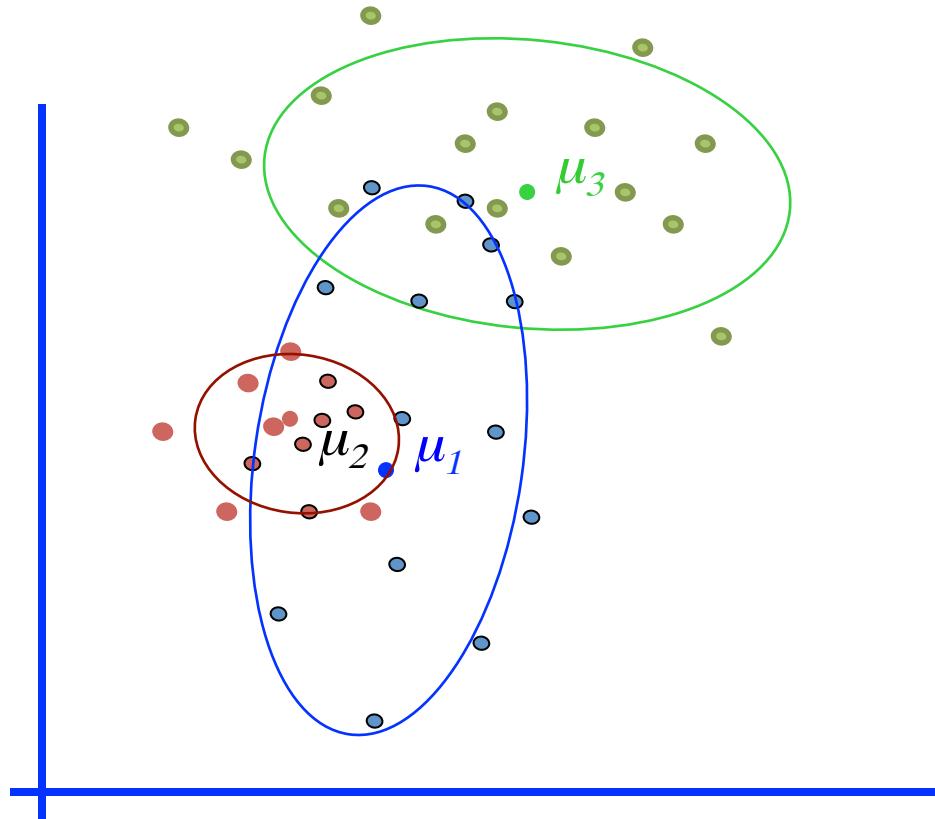
Y	X_1	X_2
??	0.1	2.1
??	0.5	-1.1
??	0.0	3.0
??	-0.1	-2.0
??	0.2	1.5
...

The General GMM assumption

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

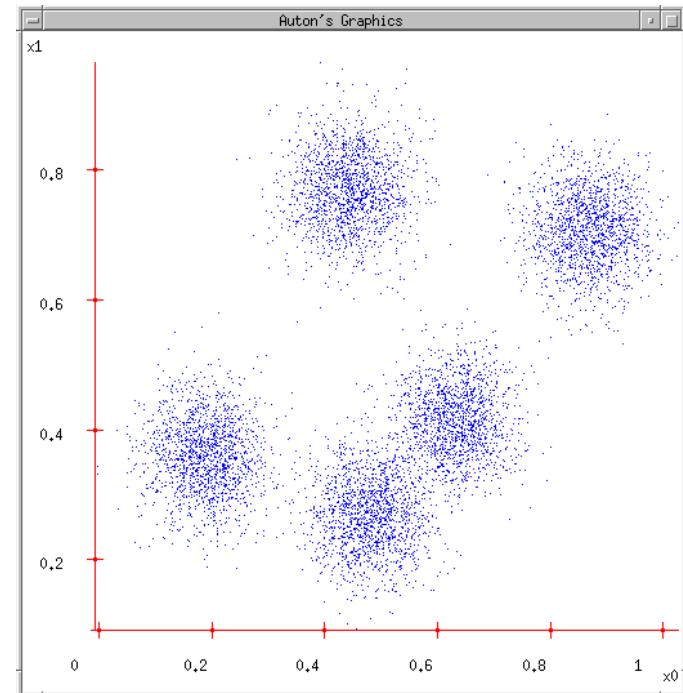
Each data point is sampled from a *generative process*:

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

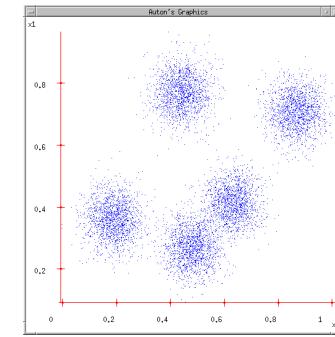


That was easy! Now, lets estimate parameters!

- MLE:
 - $\operatorname{argmax}_{\theta} \prod_j P(y^j, x^j; \theta)$
 - θ : all model parameters
 - eg, class probs, means, and variance for naïve Bayes
- But we don't know y^j !!!
- Maximize *marginal likelihood*:
 - $\operatorname{argmax}_{\theta} \prod_j P(x^j; \theta) = \operatorname{argmax} \prod_j \sum_{i=1}^k P(y^j=i, x^j; \theta)$



How do we optimize? Closed Form?

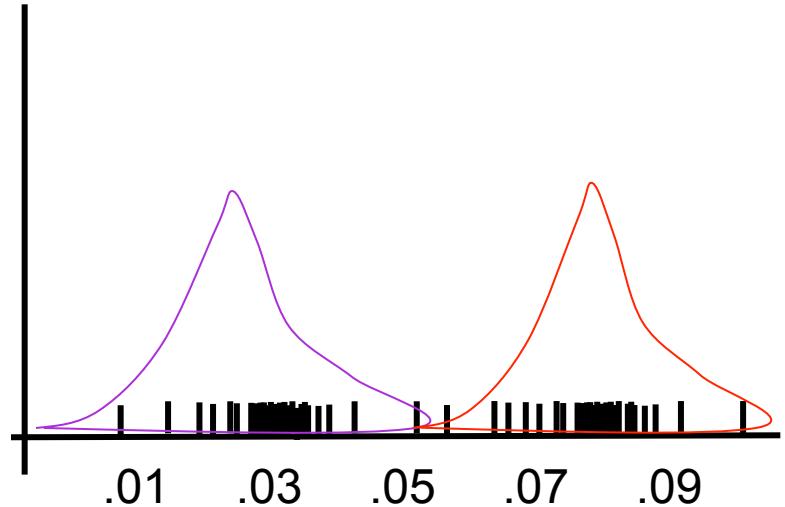


- Maximize *marginal likelihood*:
 - $\operatorname{argmax}_{\theta} \prod_j P(x^j; \theta) = \operatorname{argmax} \prod_j \sum_{i=1}^k P(y^j=i, x^j; \theta)$
- Almost always a hard problem!
 - Usually no closed form solution
 - Even when $P(X, Y; \theta)$ is convex, $P(X; \theta)$ generally isn't...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

Simple example: learn means only!

Consider:

- 1D data, m points
- Mixture of $k=2$ Gaussians
- Variances fixed to $\sigma=1$
- Dist'n over classes is uniform
- Need to estimate μ_1 and μ_2



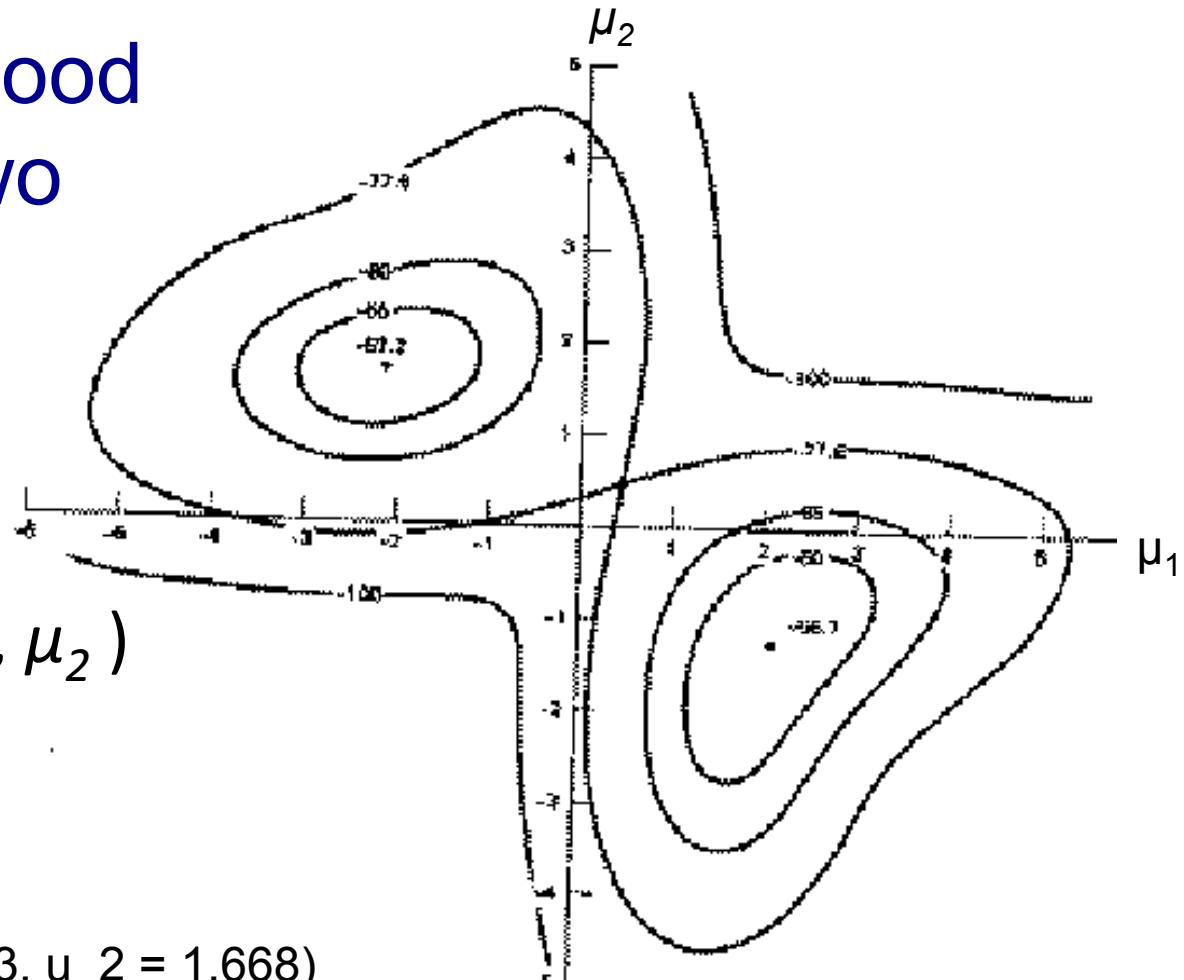
$$\prod_{j=1}^m \sum_{i=1}^k P(X = x^j, Y = i) \propto \prod_{j=1}^m \sum_{i=1}^k \exp\left(-\frac{1}{2\sigma^2}(x^j - \mu_i)^2\right)$$

Marginal Likelihood for Mixture of two Gaussians

Graph of

$\log P(x^1, x^2 \dots x^m | \mu_1, \mu_2)$

against μ_1 and μ_2



Max likelihood = ($\mu_1 = 2.13, \mu_2 = 1.668$)

Local minimum, but very close to global at ($\mu_1 = 1.668, \mu_2 = 2.13$) *

* corresponds to switching y_1 with y_2 .

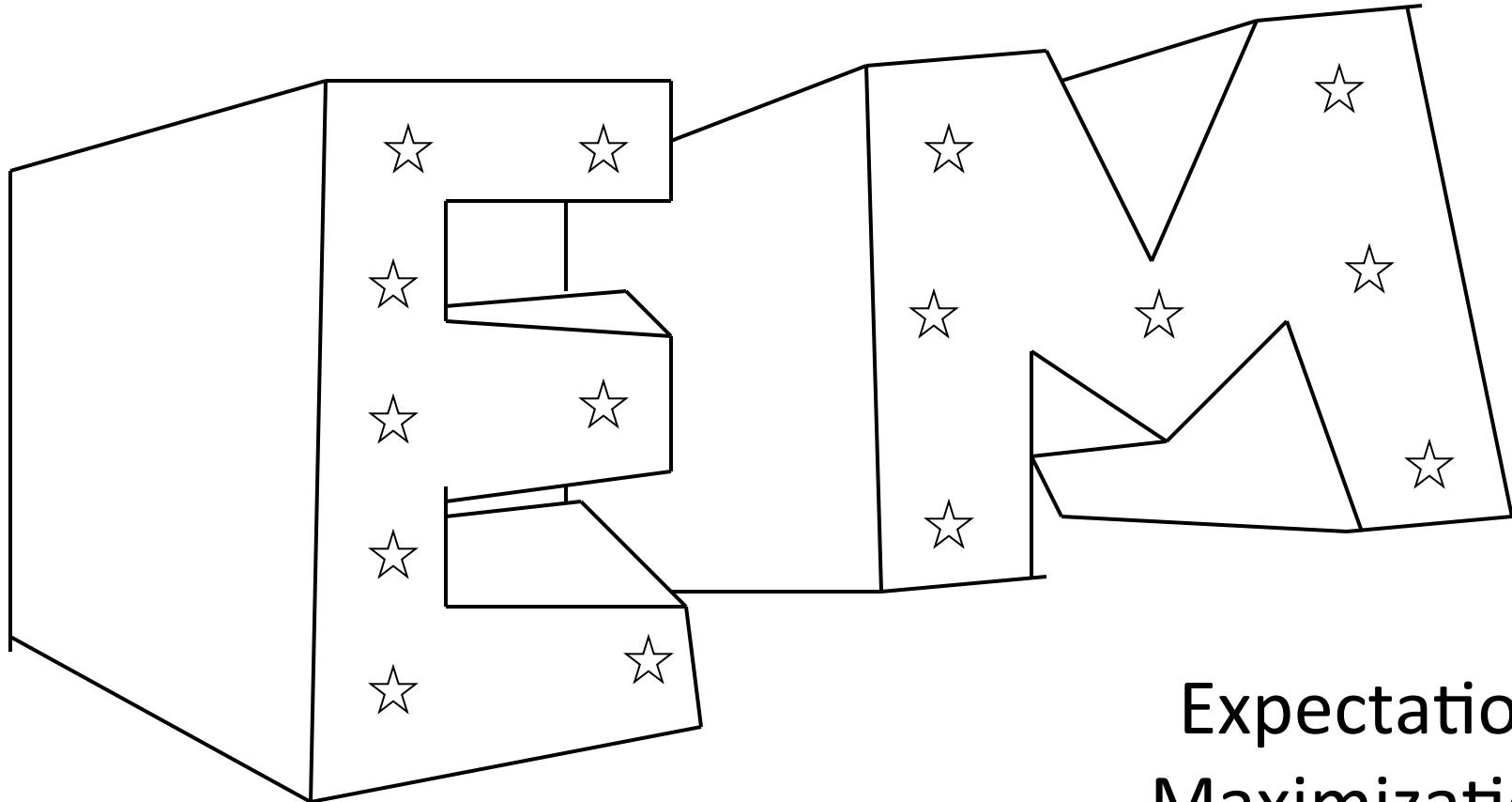
Learning general mixtures of Gaussian

$$P(X = x|Y = i) = \frac{1}{\sqrt{(2\pi)^m |\Sigma_i|}} \exp\left(-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right)$$

- Marginal likelihood, for data $\{x^j \mid j = 1..n\}$:

$$\begin{aligned} \prod_{j=1}^n P(x^j) &= \prod_{j=1}^n \sum_i P(X = x^j, Y = i) = \prod_{j=1}^n \sum_i P(X = x^j | Y = i) P(Y = i) \\ &= \prod_{j=1}^n \sum_i \frac{1}{\sqrt{(2\pi)^m |\Sigma_i|}} \exp\left(-\frac{1}{2}(x^j - \mu_i)^T \Sigma_i^{-1} (x^j - \mu_i)\right) P(Y = i) \end{aligned}$$

- Need to differentiate and solve for μ_i , Σ_i , and $P(Y=i)$ for $i=1..k$
- There will be no closed form solution, gradient is complex, lots of local optimum
- Wouldn't it be nice if there was a better way!



Expectation
Maximization

The EM Algorithm

- A clever method for maximizing marginal likelihood:
 - $\operatorname{argmax}_{\theta} \prod_j P(x^j) = \operatorname{argmax}_{\theta} \prod_j \sum_{i=1}^k P(y^j=i, x^j)$
 - A type of gradient ascent that can be 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: $\operatorname{argmax}_{\theta} \prod_j \sum_{i=1}^k P(y^j=i, x^j | \theta) = \sum_j \log \sum_{i=1}^k P(y^j=i, x^j | \theta)$

Data: $\{x^j \mid j=1 .. n\}$

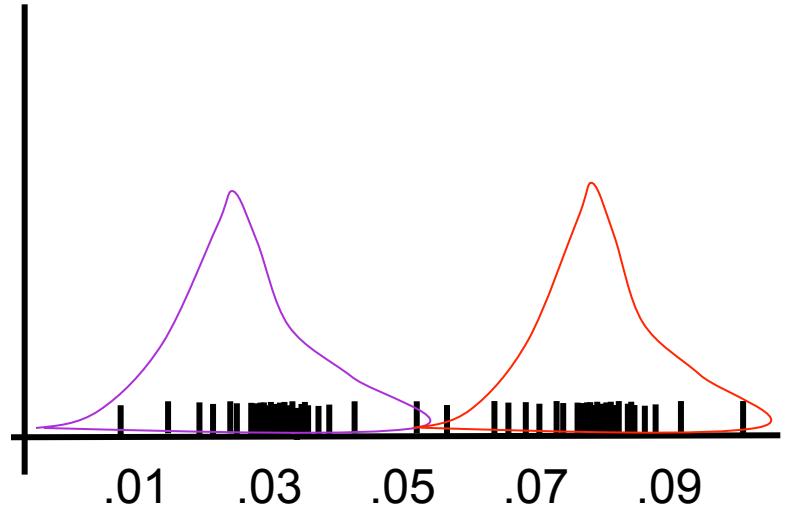
- **E-step:** Compute expectations to “fill in” missing y values according to current parameters
 - For all examples j and values i for y , compute: $P(y^j=i \mid x^j, \theta)$
- **M-step:** Re-estimate the parameters with “weighted” MLE estimates
 - Set $\theta = \operatorname{argmax}_{\theta} \sum_j \sum_{i=1}^k P(y^j=i \mid x^j, \theta) \log P(y^j=i, x^j | \theta)$

Especially useful when the E and M steps have closed form solutions!!!

Simple example: learn means only!

Consider:

- 1D data, m points
- Mixture of $k=2$ Gaussians
- Variances fixed to $\sigma=1$
- Dist'n over classes is uniform
- Need to estimate μ_1 and μ_2



$$\prod_{j=1}^n \sum_{i=1}^k P(X = x^j, Y = i) \propto \prod_{j=1}^n \sum_{i=1}^k \exp\left(-\frac{1}{2\sigma^2}(x^j - \mu_i)^2\right)$$

EM for GMMs: only learning means

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

$$\theta_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)} \}$$

E-step

Compute “expected” classes of all datapoints

$$p(y = i|x^j; \theta_t) \propto \exp\left(-\frac{1}{2\sigma^2}(x^j - \mu_i)^2\right)$$

M-step

Compute most likely new μ s given class expectations, by doing weighted ML estimates:

$$\mu_i = \frac{\sum_{j=1}^m p(y = i|x^j; \theta_t)x^j}{\sum_{j=1}^m p(y = i|x^j; \theta_t)}$$

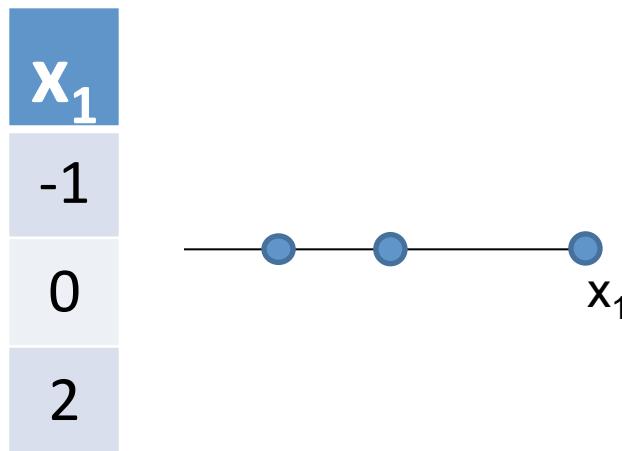
Pick K random cluster centers, $\mu_1 \dots \mu_k$
 For $t=1..T$:

- E step:

$$p(y = i|x^j; \theta_t) \propto \exp\left(-\frac{1}{2\sigma^2}(x^j - \mu_i)^2\right)$$

- M step:

$$\mu_i = \frac{\sum_{j=1}^m p(y = i|x^j; \theta_t)x^j}{\sum_{j=1}^m p(y = i|x^j; \theta_t)}$$



Initialization, random means and $\sigma=1$:

- $\mu_1 = -1, \mu_2 = 0$

$t=0$:

- $P(y=1|x^1) \propto \exp(-0.5 \times (-1+1)^2) = 1$
- $P(y=2|x^1) \propto \exp(-0.5 \times (-1-0)^2) = 0.6$
 - $P(y=1|x^1) = 0.63, P(y=2|x^1) = 0.37$
- $P(y=1|x^2) \propto \exp(-0.5 \times (0+1)^2) = 0.6$
- $P(y=2|x^2) \propto \exp(-0.5 \times (0-0)^2) = 1$
 - $P(y=1|x^2) = 0.37, P(y=2|x^2) = 0.63$
- $P(y=1|x^3) \propto \exp(-0.5 \times (2+1)^2) = 0.07$
- $P(y=2|x^3) \propto \exp(-0.5 \times (2-0)^2) = 0.93$
 - $P(y=1|x^3) = 0.01, P(y=2|x^3) = 0.93$
- $\mu^1 = (0.63 \times -1 + 0.37 \times 0 + 0.07 \times 2) / (0.63 + 0.37 + 0.07) = -0.45$
- $\mu^2 = (0.37 \times -1 + 0.67 \times 0 + 0.93 \times 2) / (0.37 + 0.67 + 0.93) = 0.75$

$t=1$:

- learning continues, when do we stop?

E.M. for General GMMs

Iterate: On the t 'th iteration let our estimates be, for y with k classes

$$\theta_t = \{ \mu_1 \dots \mu_k, \Sigma_1 \dots \Sigma_k, p_1, \dots, p_k \}$$

E-step

Compute “expected” classes of all datapoints for each class

$$P(y = i|x^j; \theta_t) \propto \frac{1}{\sqrt{(2\pi)^m |\Sigma_i|}} \exp\left(-\frac{1}{2}(x^j - \mu_i)^T \Sigma_i^{-1} (x^j - \mu_i)\right) p_i$$

Evaluate a Gaussian at x^j

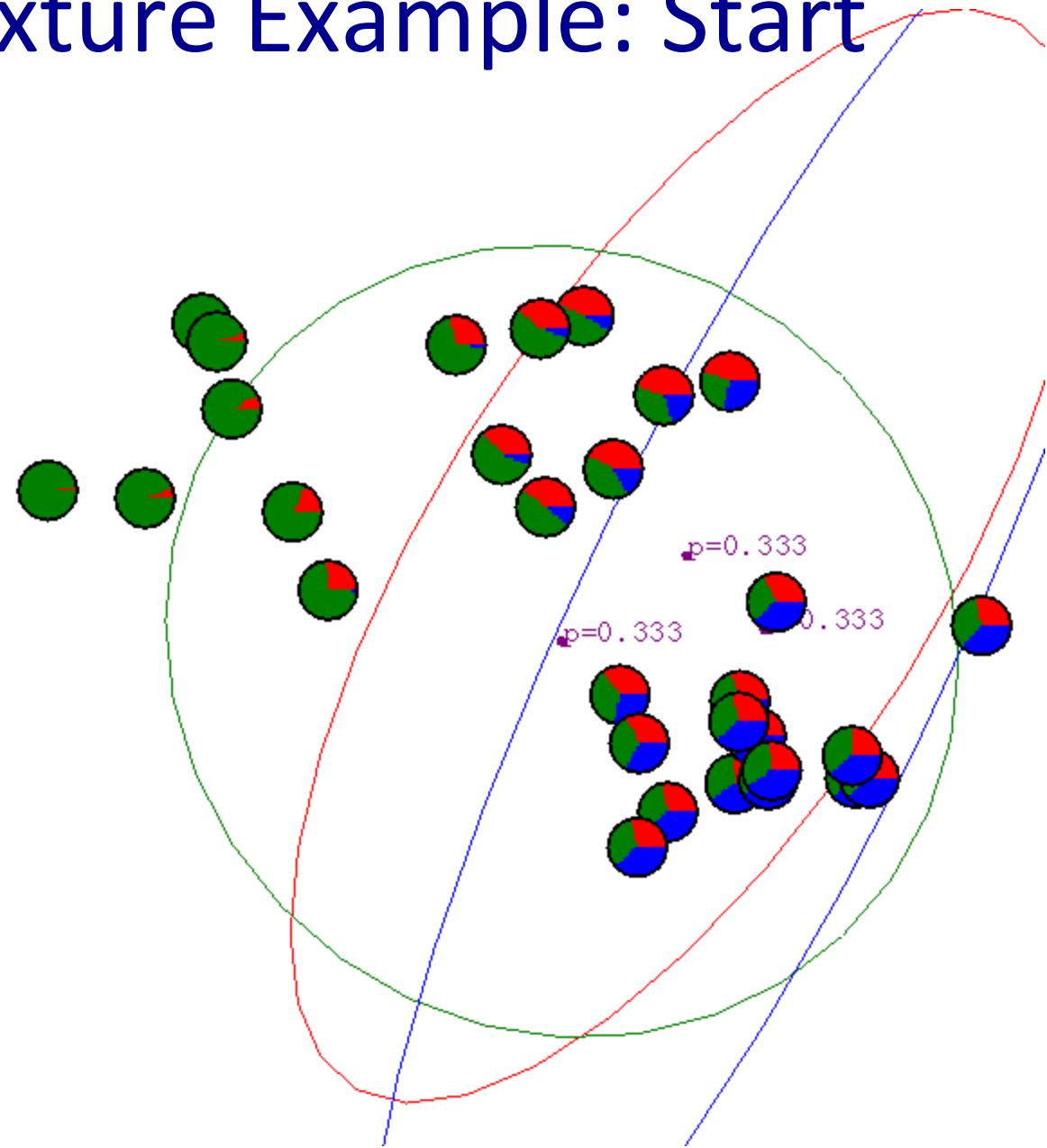
M-step

Compute weighted MLE for μ and Σ given expected classes above

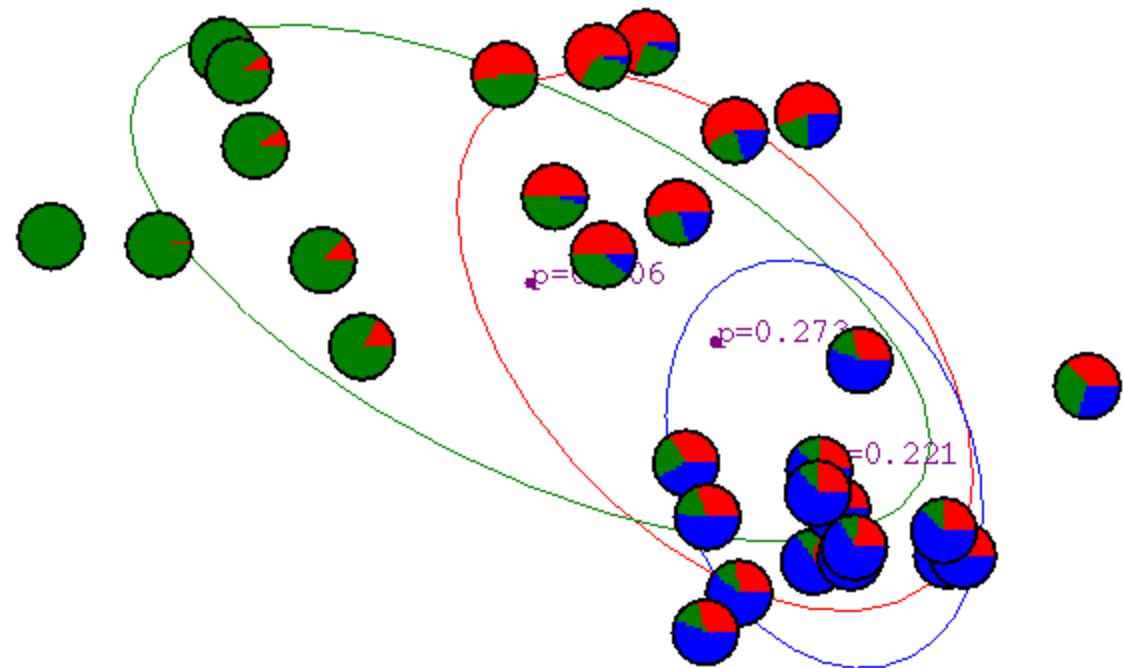
$$\mu_i = \frac{\sum_{j=1}^m p(y = i|x^j; \theta_t)x^j}{\sum_{j=1}^m p(y = i|x^j; \theta_t)} \quad \Sigma_i = \frac{\sum_{j=1}^m p(y = i|x^j; \theta_t)(x^j - \mu_i)(x^j - \mu_i)^T}{\sum_{j=1}^m p(y = i|x^j; \theta_t)}$$

$$p_i = \frac{1}{m} \sum_{j=1}^m p(y = i|x^j; \theta_t)$$

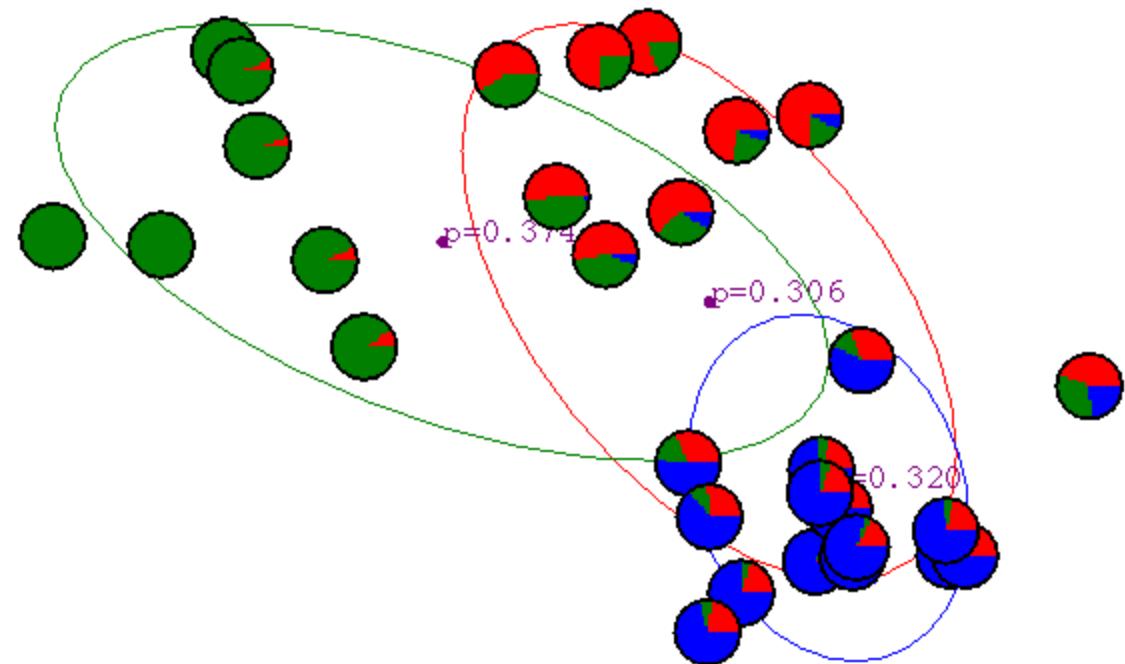
Gaussian Mixture Example: Start



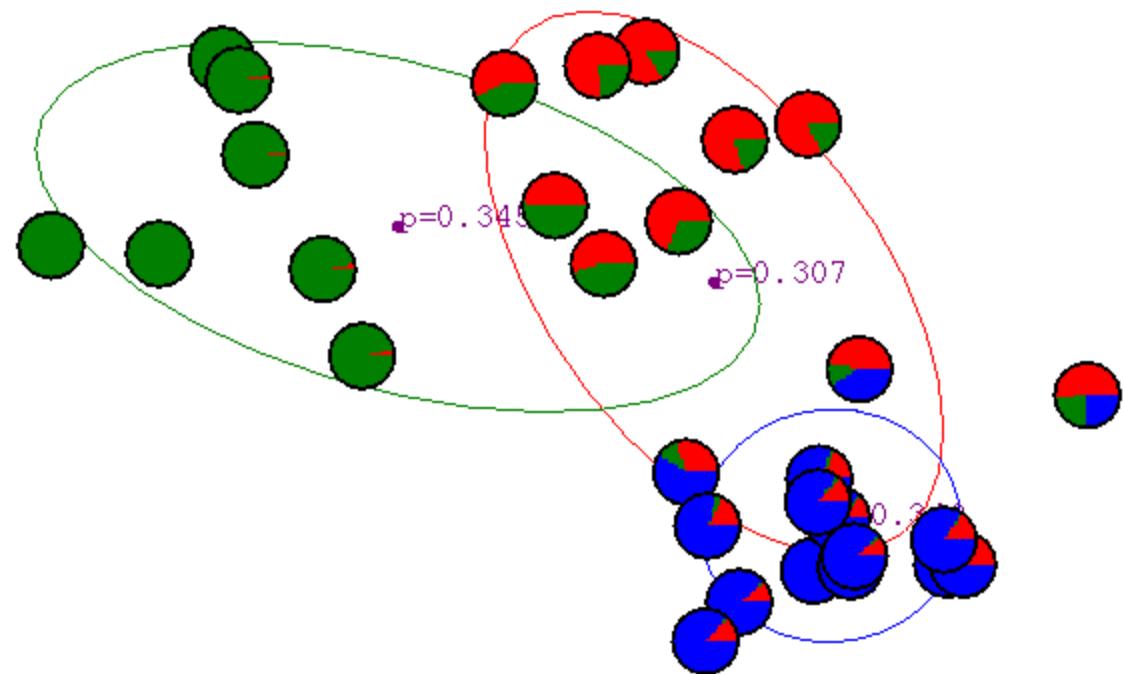
After first iteration



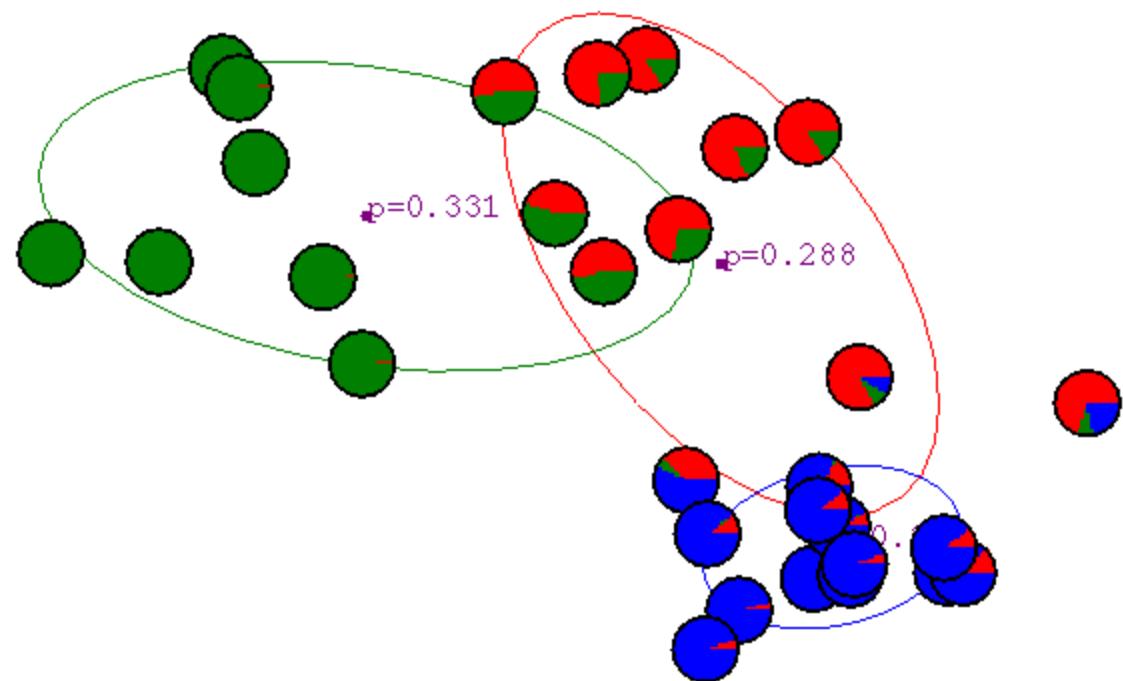
After 2nd iteration



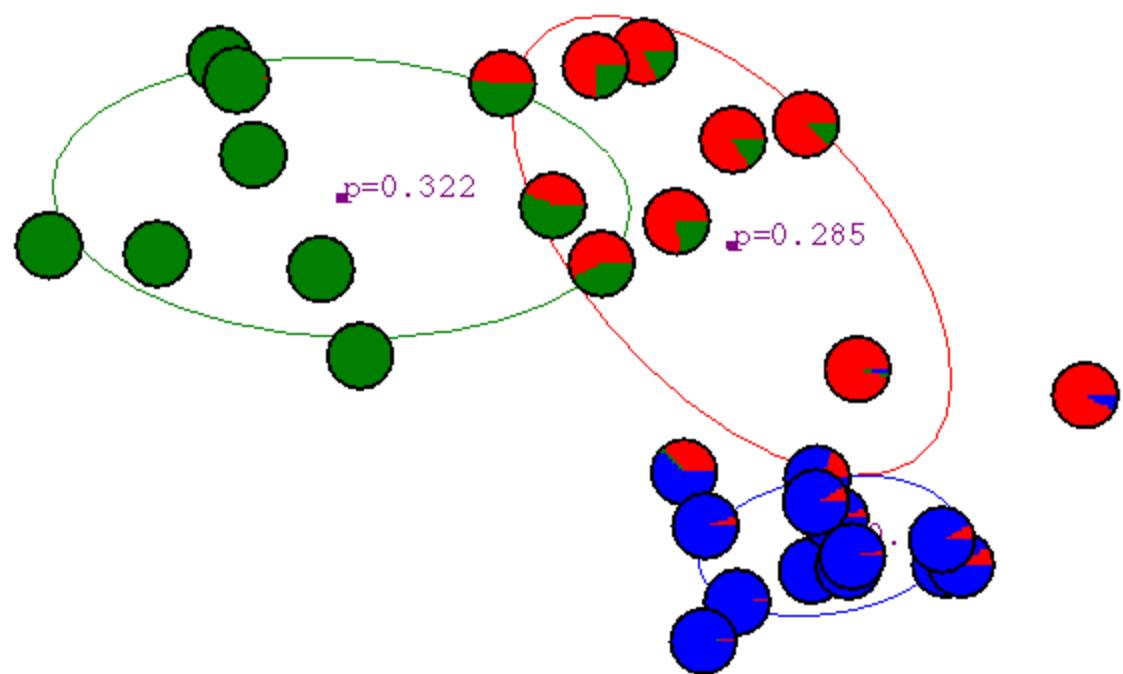
After 3rd iteration



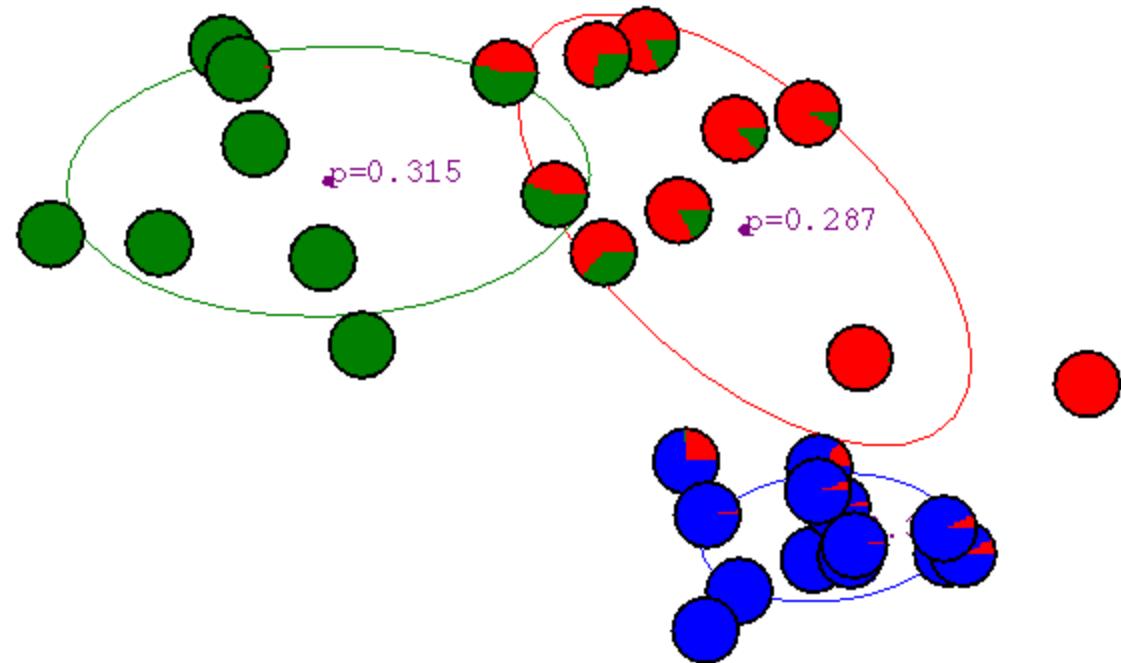
After 4th iteration



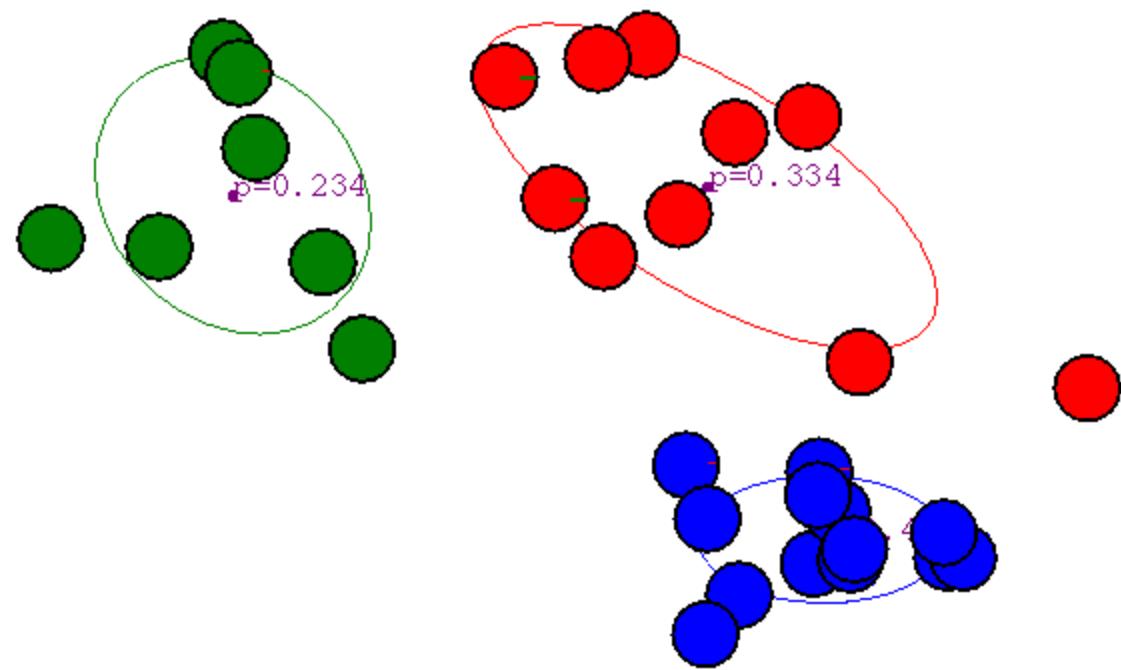
After 5th iteration



After 6th iteration



After 20th iteration



What if we do hard assignments, and learn means only?

E-step / Compute cluster assignment

Compute “expected” classes → set most likely class

$$p(y = i|x^j; \theta_t) = \exp\left(-\frac{1}{2\sigma^2}\|x^j - \mu_i\|_2^2\right) \quad \rightarrow \quad a^i = \arg \min_j dist(x^i, c^j)$$

$dist(x, x') = \|x - x'\|_2^2$

M-step / Recompute cluster mean

Compute most likely new μ s → averages over hard assignments

$$\mu_i = \frac{\sum_{j=1}^m p(y = i|x^j; \theta_t)x^j}{\sum_{j=1}^m p(y = i|x^j; \theta_t)} \quad \rightarrow \quad c^i = \frac{1}{|\{j|a^j = i\}|} \sum_{j|a^j=i} x^j$$

With hard assignments and unit variance, EM is equivalent to k-means clustering algorithm!!!

EM: Two Easy Steps

Objective: $\operatorname{argmax}_{\theta} \prod_j \sum_{i=1}^k P(y^j=i, x^j | \theta) = \sum_j \log \sum_{i=1}^k P(y^j=i, x^j | \theta)$

Data: $\{x^j \mid j=1 .. n\}$

- **E-step:** Compute expectations to “fill in” missing y values according to current parameters
 - For all examples j and values i for y , compute: $P(y^j=i \mid x^j, \theta)$
- **M-step:** Re-estimate the parameters with “weighted” MLE estimates
 - Set $\theta = \operatorname{argmax}_{\theta} \sum_j \sum_{i=1}^k P(y^j=i \mid x^j, \theta) \log P(y^j=i, x^j | \theta)$

Especially useful when the E and M steps have closed form solutions!!!

Implementation Detail: Working with Log Probabilities

$$P(c|x_1, \dots, x_n) = \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

Implementation Detail: Working with Log Probabilities

$$P(c|x_1, \dots, x_n) = \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

$$\log P(c|x_1, \dots, x_n) = \log \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

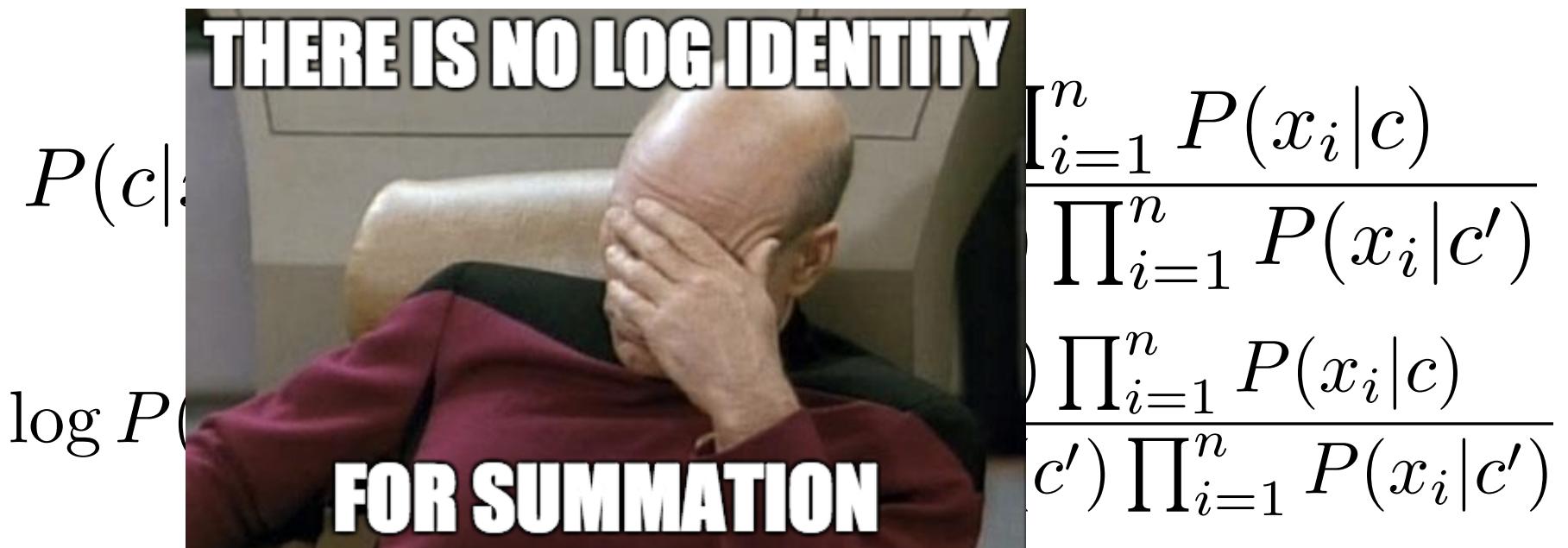
Implementation Detail: Working with Log Probabilities

$$P(c|x_1, \dots, x_n) = \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

$$\log P(c|x_1, \dots, x_n) = \log \frac{P(c) \prod_{i=1}^n P(x_i|c)}{\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c')}$$

$$= \log P(c) + \sum_{i=1}^n \log P(x_i|c) - \log \left[\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c') \right]$$

Implementation Detail: Working with Log Probabilities



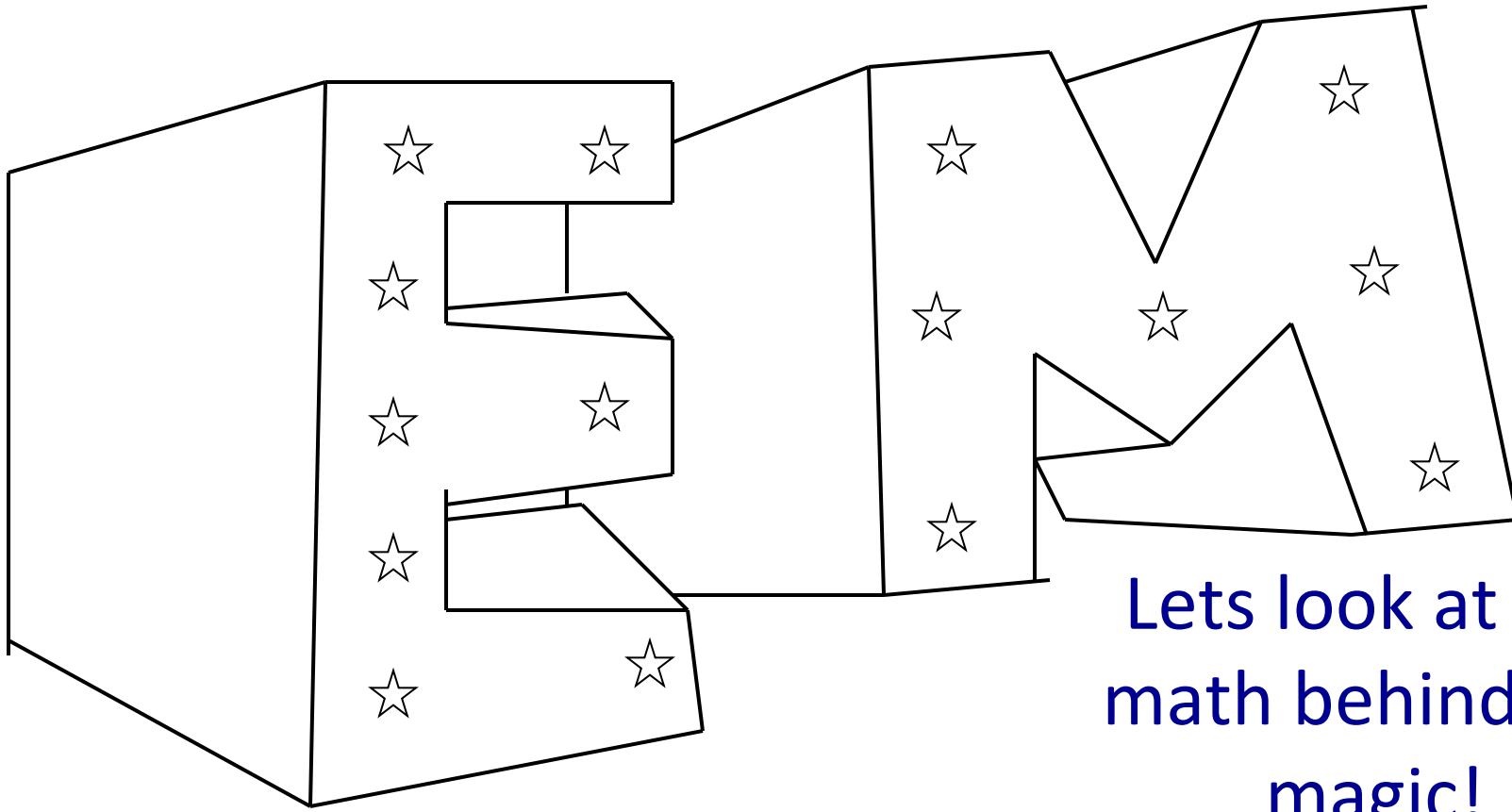
$$= \log P(c) + \sum_{i=1}^n \log P(x_i|c) - \log \left[\sum_{c'} P(c') \prod_{i=1}^n P(x_i|c') \right]$$

Log Exp Sum Trick: motivation

- We have: a bunch of log probabilities.
 - $\log(p_1), \log(p_2), \log(p_3), \dots \log(p_n)$
- We want: $\log(p_1 + p_2 + p_3 + \dots p_n)$
- We could convert back from log space,
sum then take the log.
 - If the probabilities are very small, this will
result in floating point underflow

Log Exp Sum Trick:

$$\log\left[\sum_i \exp(x_i)\right] = x_{max} + \log\left[\sum_i \exp(x_i - x_{max})\right]$$



Lets look at the
math behind the
magic!

We will argue that EM:

- Optimizes a bound on the likelihood
- Is a type of coordinate ascent
- Is guaranteed to converge to a (often local) optima

The general learning problem with missing data

- Marginal likelihood: \mathbf{x} is observed, \mathbf{z} is missing:

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

- Objective: Find $\operatorname{argmax}_{\theta} \ell(\theta : \mathcal{D})$

A Key Computation: E-step

- \mathbf{x} is observed, \mathbf{z} is missing
- Compute probability of missing data given current choice of θ
 - $Q(\mathbf{z} | \mathbf{x}_j)$ for each \mathbf{x}_j
 - e.g., probability computed during classification step
 - corresponds to “classification step” in K-means

$$Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) = P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$$

Jensen's inequality

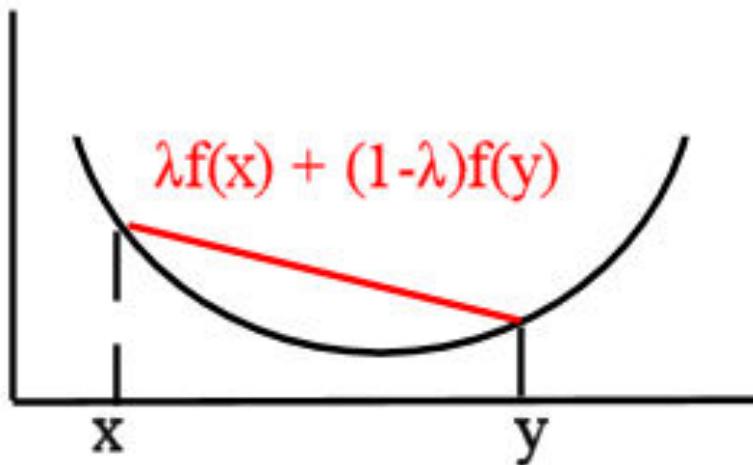
$$\ell(\theta : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{z} | \mathbf{x}_j) P(\mathbf{x}_j | \theta)$$

- **Theorem:**

- $-\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

- e.g., **Binary case for convex function f:**

$$\lambda f(x) + (1 - \lambda)f(y) \geq f(\lambda x + (1 - \lambda)y)$$



- actually, holds for any concave (convex) function applied to an expectation!

Applying Jensen's inequality

- Use: $\log \sum_z P(z) f(z) \geq \sum_z P(z) \log f(z)$

$$\ell(\theta^{(t)} : \mathcal{D}) = \sum_{j=1}^m \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)}$$

$$\geq \sum_{j=1}^m \sum_z Q^{(t+1)}(z | x_j) \log \left(\frac{p(z, x_j | \theta^{(t)})}{Q^{(t+1)}(z | x_j)} \right)$$

$$= \sum_{j=1}^m \sum_z Q^{(t+1)}(z | x_j) \log \left(p(z, x_j | \theta^{(t)}) \right) - \sum_{j=1}^m \sum_z Q^{(t+1)}(z | x_j) \log \left(Q^{(t+1)}(z | x_j) \right)$$

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

The M-step

Lower bound:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + m.H(Q^{(t+1)})$$

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)$$

- We are optimizing a lower bound!
- Use expected counts to do weighted learning:
 - If learning requires $\text{Count}(\mathbf{x}, \mathbf{z})$
 - Use $E_{Q^{(t+1)}}[\text{Count}(\mathbf{x}, \mathbf{z})]$
 - *Looks a bit like boosting!!!*

Convergence of EM

- Define: potential function $F(\theta, Q)$:

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

- lower bound from Jensen's inequality
- EM is coordinate ascent on F !
 - Thus, maximizes lower bound on marginal log likelihood

M-step can't decrease $F(\theta, Q)$: by definition!

$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta)$$

- We are maximizing F directly, by ignoring a constant!

$$F(\theta, Q^{(t+1)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t+1)})$$

E-step: more work to show that $F(\theta, Q)$ doesn't decrease

- KL-divergence: measures distance between distributions

$$KL(Q||P) = \sum_z Q(z) \log \frac{Q(z)}{P(z)}$$

- KL=zero if and only if $Q=P$

E-step also doesn't decrease F: Step 1

- Fix θ to $\theta^{(t)}$, take a max over Q:

$$\begin{aligned}\ell(\theta^{(t)} : \mathcal{D}) &\geq F(\theta^{(t)}, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)} \\ &= \sum_{j=1}^m \sum_z Q(z \mid x_j) \log \left(\frac{P(z \mid x_j, \theta^{(t)}) P(x_j \mid \theta^{(t)})}{Q(z \mid x_j)} \right) \\ &= \sum_{j=1}^m \sum_z Q(z \mid x_j) \log \left(P(x_j \mid \theta^{(t)}) \right) - \sum_{j=1}^m \sum_z Q(z \mid x_j) \log \left(\frac{Q(z \mid x_j)}{P(z \mid x_j, \theta^{(t)})} \right) \\ &= \ell(\theta^{(t)} : \mathcal{D}) - \sum_{j=1}^m KL \left(Q(\mathbf{z} \mid \mathbf{x}_j) \middle\| P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)}) \right)\end{aligned}$$

E-step also doesn't decrease F: Step 2

- Fixing θ to $\theta^{(t)}$:

$$\begin{aligned}\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) &= \ell(\theta^{(t)} : \mathcal{D}) + \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} | \mathbf{x}_j)} \\ &= \ell(\theta^{(t)} : \mathcal{D}) - \sum_{j=1}^m KL(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}))\end{aligned}$$

- Now, the max over Q yields:

- $Q(\mathbf{z} | \mathbf{x}_j) \leftarrow P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)})$
- Why? The likelihood term is a constant; the KL term is zero iff the arguments are the same distribution!!
- So, the E-step is actually a maximization / tightening of the bound. It ensures that:

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

EM is coordinate ascent

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q) = \sum_{j=1}^m \sum_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j | \theta)}{Q(\mathbf{z} | \mathbf{x}_j)}$$

- **M-step:** Fix Q , maximize F over θ (a lower bound on $\ell(\theta : \mathcal{D})$):

$$\ell(\theta : \mathcal{D}) \geq F(\theta, Q^{(t)}) = \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta) + m.H(Q^{(t)})$$

- **E-step:** Fix θ , maximize F over Q :

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^m KL \left(Q(\mathbf{z} | \mathbf{x}_j) || P(\mathbf{z} | \mathbf{x}_j, \theta^{(t)}) \right)$$

- “Realigns” F with likelihood:

$$F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})$$

What you should know

- K-means for clustering:
 - algorithm
 - converges because it's coordinate ascent
- Know what agglomerative clustering is
- EM for mixture of Gaussians:
 - How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it DOES
- EM is coordinate ascent
- General case for EM

Acknowledgements

- K-means & Gaussian mixture models presentation contains material from excellent tutorial by Andrew Moore:
 - <http://www.autonlab.org/tutorials/>
- K-means Applet:
 - http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/AppletKM.html
- Gaussian mixture models Applet:
 - <http://www.neurosci.aist.go.jp/%7Ekaho/MixtureEM.html>