



ISTA 421 + INFO 521

Introduction to Machine Learning

Lecture 23: Clustering

**K-Means,
Kernelized K-Means
Mixture Models**

Clay Morrison

claytonm@email.arizona.edu

Gould-Simpson 819

Phone 621-6609

17 November 2015

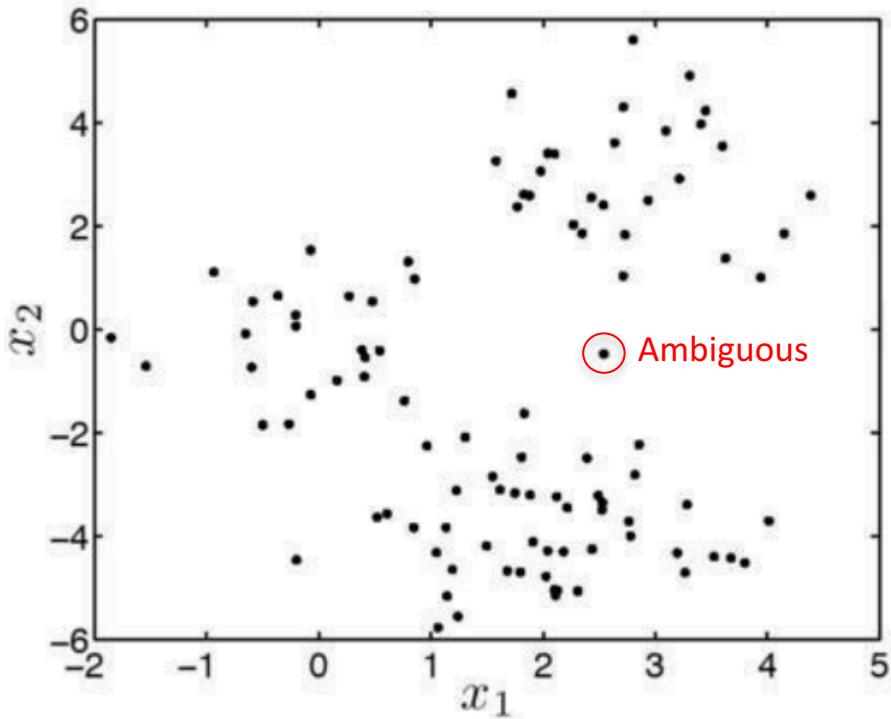
- Minimizing a Loss function
 - Linear model
 - Linear least mean squares
 - Maximum Likelihood
 - Probabilistic model of uncertainty (noise, error)
 - Maximize the likelihood w.r.t. parameters
 - Linear model with additive Gaussian noise
 - Bayesian Approach
 - Treat parameters as random variables
 - Use Bayes Theorem to combine likelihood & prior to find posterior distribution
 - Estimation Techniques (often used in Bayesian approaches)
 - Gradient methods (Widrow-Hoff (1st), Newton-Raphson (2nd))
 - Laplace Approximation (estimate posterior with Gaussian)
 - Monte Carlo estimation of expectation; Metropolis-Hastings
 - Classification (& Regression)
 - Clustering
 - Projection
- Approaches to Avoiding Over-fitting
i.e., how to max **generalization** performance
- Regularization (controlling model complexity)
- Cross Validation (estimating the gen error)
- Marginal Likelihood model selection
- predicting output
- Main algorithmic families of Machine Learning

Clustering

- **Unsupervised learning:** only provided with set of objects \times
- Goal of **cluster analysis:** create grouping of objects where objects within group are similar and objects in different groups are not similar (or as similar).
- Examples
 - **Customer preference:** lots of data about purchases, used as evidence for personal preferences. Define measure of similarity between customers based on purchasing history: within group similar shopping patterns. Recommend items based on customer similarity. Also cluster items based on customers they were purchased by (items 1 and 2 both bought by customers A, D, E and G).
 - **Gene function prediction:** categorize genes into functional classes based on patterns of mutual interaction in mRNA microarray data. Functions of known genes in cluster can be used as predicted function of unknown genes in same cluster.

Example

How many groups?



Many ways of defining similarity in terms of (inverse) distance:

For real valued data...

Euclidean distance:

$$(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)$$

First approach: characterize clusters by centroid;
members of clusters are closest to cluster centroid.

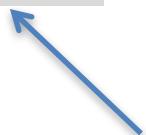
K-means Clustering

- Mean point of k^{th} cluster: μ_k
- Binary indicator function: z_{nk}
 - 1 if object n is assigned to cluster k , 0 otherwise
 - Each object assigned to one and only one cluster

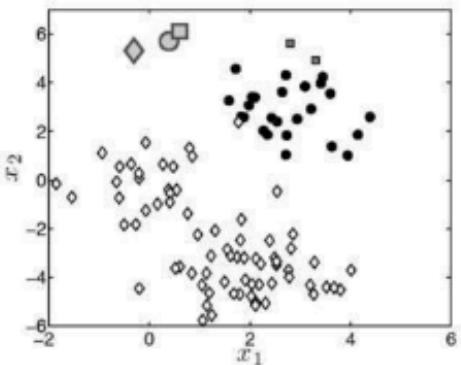
Algorithm:

- (0) Choose K (total number of clusters) and initial random cluster means μ_1, \dots, μ_K
- (1) For each data object, find closest cluster k mean and set $z_{nk} = 1$ and $z_{nj} = 0$ for all $j \neq k$.
- (2) If all the assignments (z_{nk}) are unchanged from the previous iteration, stop.
- (3) Update each μ_k
- (4) Goto 1

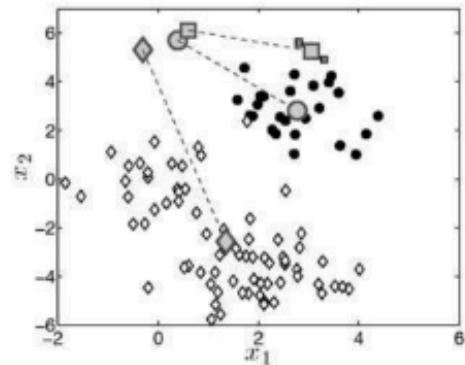
$$\mu_k = \frac{\sum_n z_{nk} \mathbf{x}_n}{\sum_n z_{nk}}$$



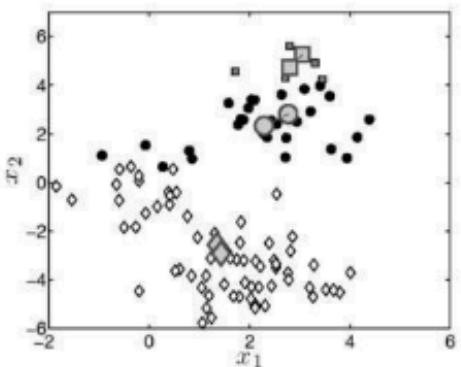
K-means Example



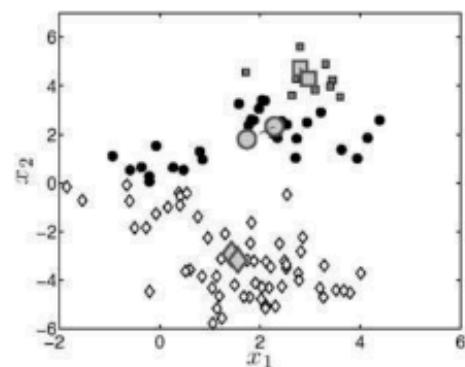
(a) Data and initial random means.
Means are depicted by large symbols.
Each data object is given the symbol of
its closest mean



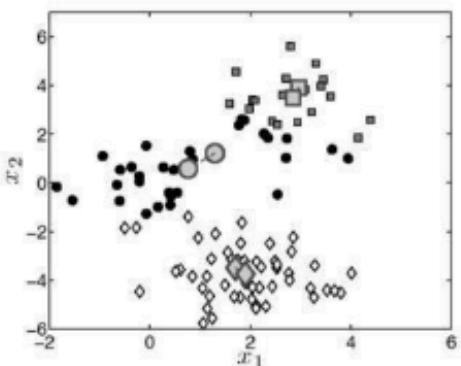
(b) Means updated according to assigned
objects



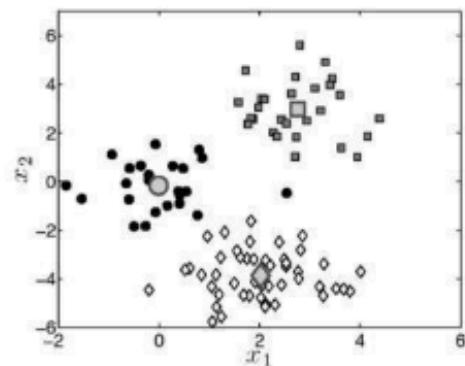
(c) Objects re-assigned to new means and
means updated again



(d) Means updated after three iterations



(e) Means updated after five iterations



(f) Means updated after eight iterations.
Algorithm has converged

Also, simple Java applet

http://www.math.le.ac.uk/people/ag153/homepage/KmeansKmedoids/Kmeans_Kmedoids.html

K -means Clustering

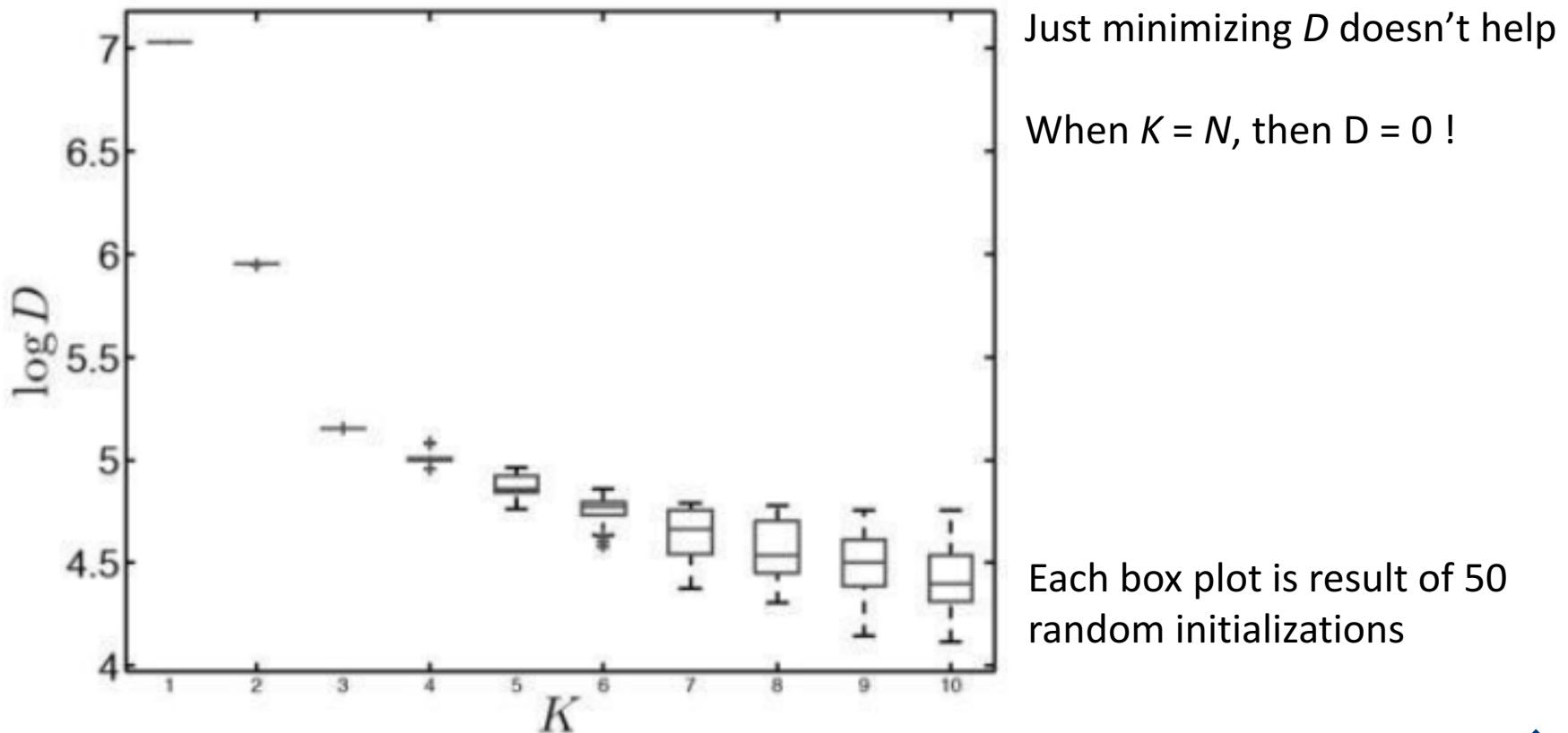
- Finds **local** minimum of total distance of all points to their cluster centers:

$$D = \sum_{n=1}^N \sum_{k=1}^K z_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Whether finds **global** minimum depends on initial cluster means – can get caught in local minima
- Guaranteed to find global minimum only if evaluate every possible assignment of all N points to K clusters – intractable
- Typically just run multiple times, select final cluster means with lowest total distance D

K -means Clustering

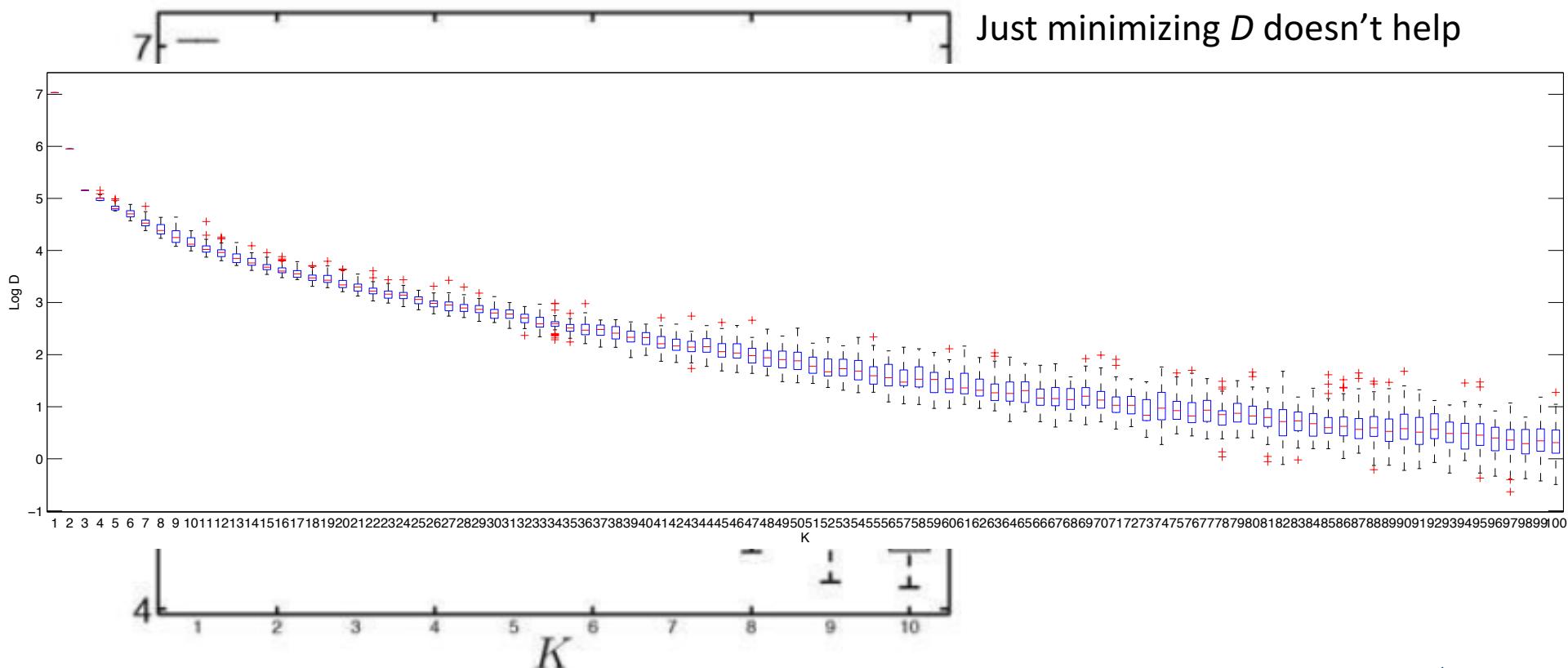
- Choosing the number of clusters: K
 - A common problem in cluster analysis



Since cannot optimize directly, an often used technique for choosing K is to see how use of the clusters works on **other** performance tasks

K -means Clustering

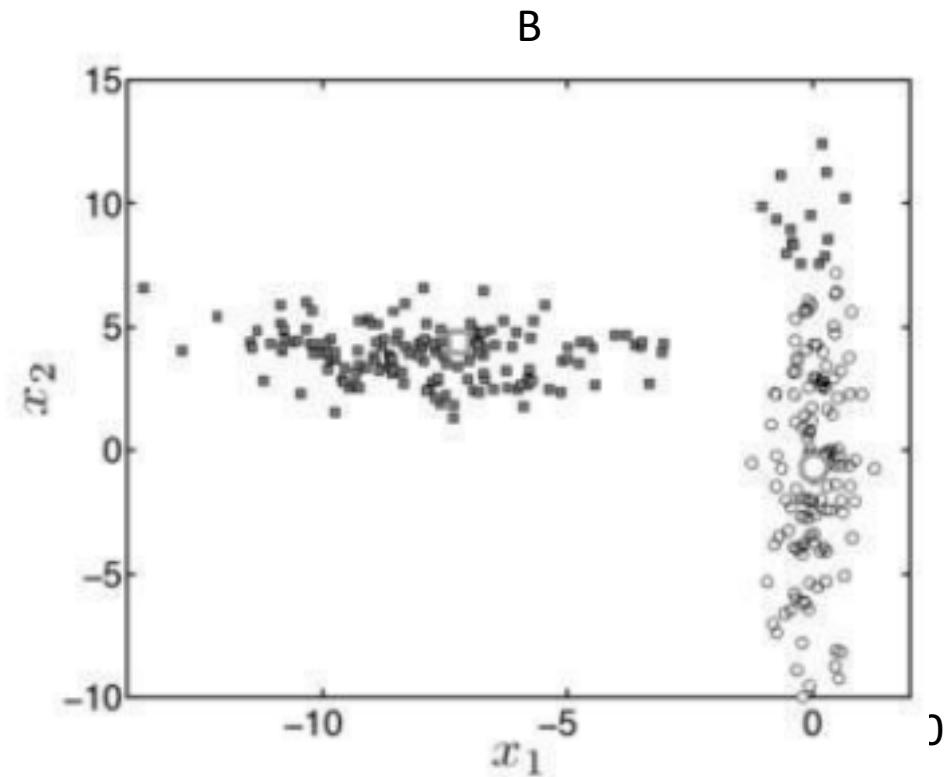
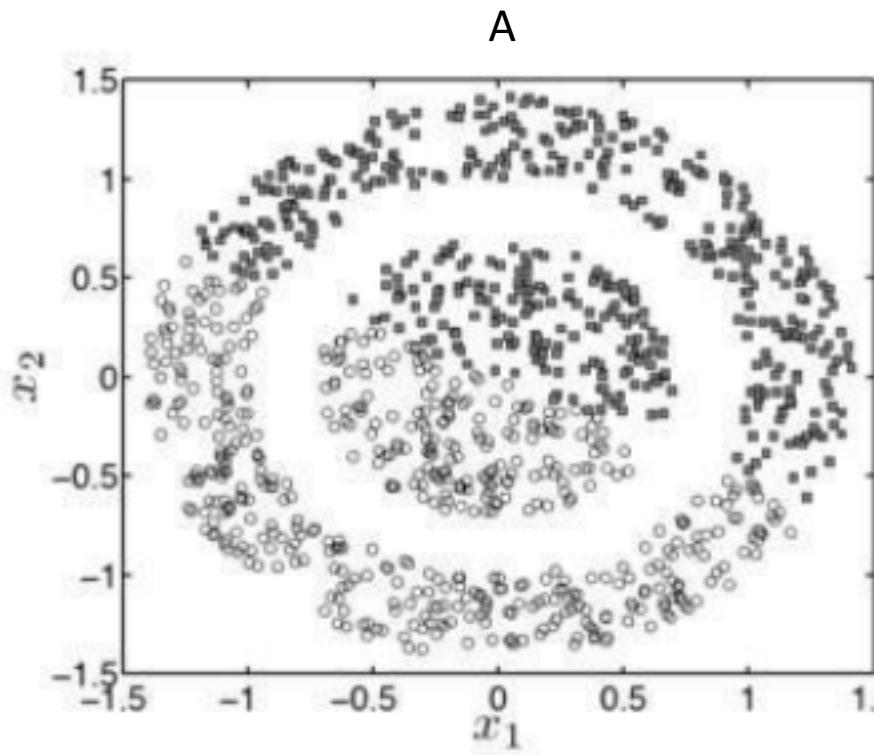
- Choosing the number of clusters: K
 - A common problem in cluster analysis



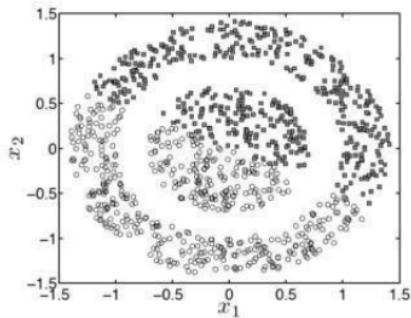
Since cannot optimize directly, an often used technique for choosing K is to see how use of the clusters works on other performance tasks

Where K-means fails

- Objects in true clusters do not necessarily conform to current distance-based similarity



Kernelized K-means



- A kernelized K-means can handle case A
- Derive kernelizable form of distance:

$$d_{nk} = (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

$$N_k = \sum_{n=1}^N z_{nk}$$

$$d_{nk} = \left(\mathbf{x}_n - \frac{1}{N_k} \sum_{m=1}^N z_{mk} \mathbf{x}_m \right)^\top \left(\mathbf{x}_n - \frac{1}{N_k} \sum_{r=1}^N z_{rk} \mathbf{x}_r \right)$$

Note

NOTE: the means $\boldsymbol{\mu}_k$ doesn't appear here!

Is kernelizable!

$$d_{nk} = K(\mathbf{x}_n, \mathbf{x}_n) - \frac{2}{N_k} \sum_{m=1}^N z_{mk} K(\mathbf{x}_n, \mathbf{x}_m) + \frac{1}{N_k^2} \sum_{m=1}^N \sum_{r=1}^N z_{mk} z_{rk} K(\mathbf{x}_m, \mathbf{x}_r)$$

However, computing $\boldsymbol{\mu}_k$ directly with a transformation...

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N z_{nk} \mathbf{x}_n}{\sum_{n=1}^N z_{nk}}$$

→

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N z_{nk} \phi(\mathbf{x}_n)}{\sum_{n=1}^N z_{nk}}$$

Not kernelizable

And when transformation is not explicitly computable, can't even compute the centroid!

Kernelized K-means

- Since want to avoid directly needing to compute the means (centroids), need to augment the algorithm:

$$d_{nk} = K(\mathbf{x}_n, \mathbf{x}_n) - \frac{2}{N_k} \sum_{m=1}^N z_{mk} K(\mathbf{x}_n, \mathbf{x}_m) + \frac{1}{N_k^2} \sum_{m=1}^N \sum_{r=1}^N z_{mk} z_{rk} K(\mathbf{x}_m, \mathbf{x}_r)$$

- (1) Randomly initialize z_{nk} for each n ← Initializing by object-cluster rather than cluster means!
- (2) Compute d_{n1}, \dots, d_{nk} for each object n ← Only requires pairwise inner products (kernelized)
- (3) Assign each object to the cluster k with the lowest d_{nk}
 - This determines the new z_{nk} for each n
- (4) If assignments have changed, goto step 2, otherwise stop.

Two variants of initialization step:

- (A) Run regular K-means to convergence, then use z_{nk} as seed to Kernelized version
- (B) Assign $N - K + 1$ objects to cluster 1 and remaining $K - 1$ objects to separate clusters.

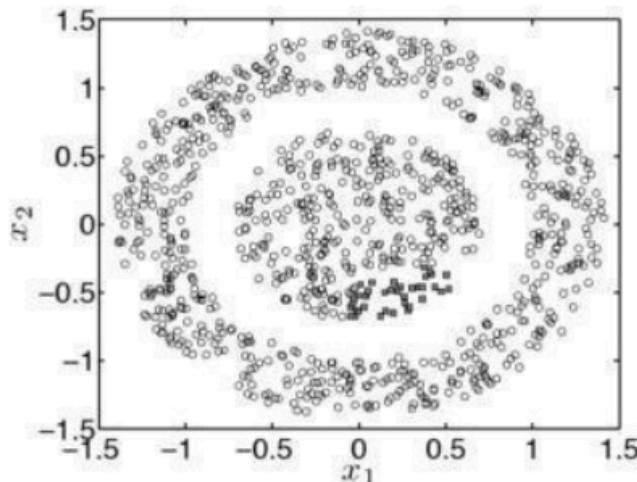
Kernel K-Means for case A

Initialization:

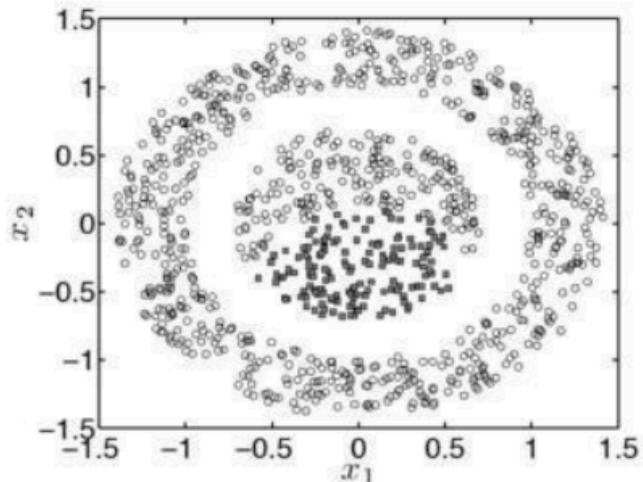
Assigned 1 point to “squares”
and
the remaining N-1 to “circles”

Gaussian Kernel: $\gamma=1$

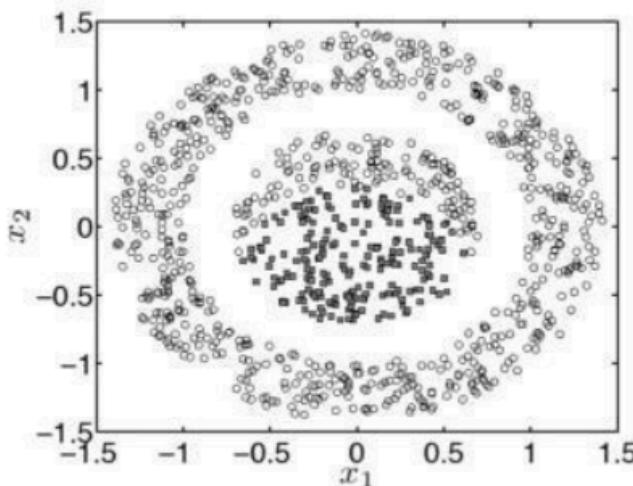
$$\exp \left\{ -\gamma(\mathbf{x}_n - \mathbf{x}_m)^T (\mathbf{x}_n - \mathbf{x}_m) \right\}$$



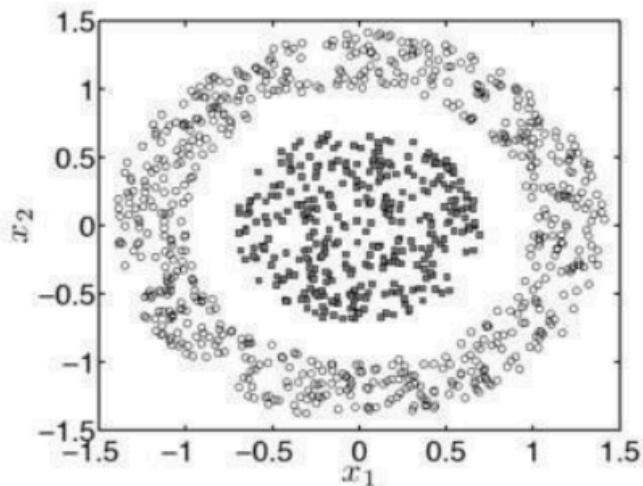
(a) Kernel K-means after one iteration



(b) After five iterations



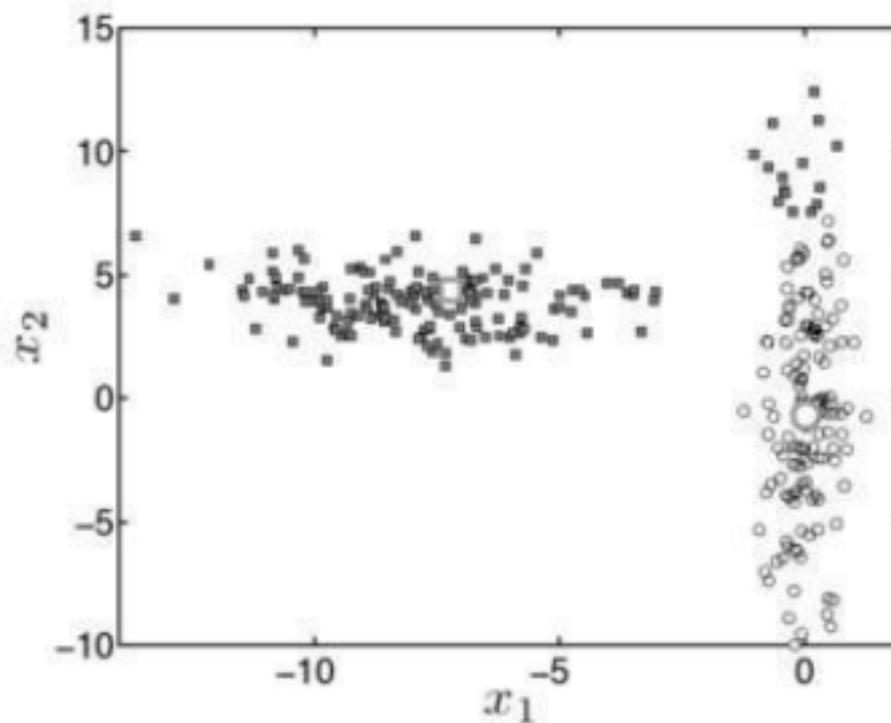
(c) After 10 iterations



(d) At convergence (30 iterations)

Mixture Models

- Some similarities to K -means, but much richer representations of the data (rather than points / centroids)



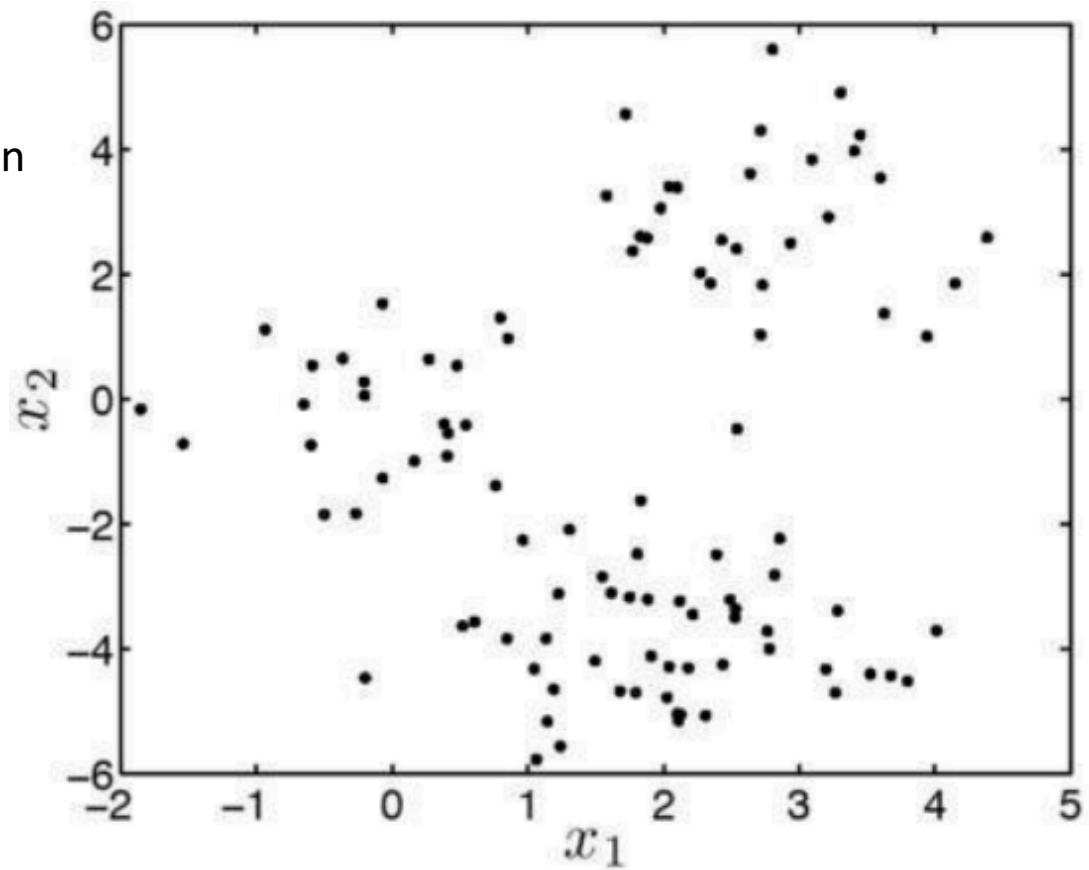
Centroids model of clusters is too simple to capture the structure here

The Generative Picture (again)

- How could we *generate* this data?

For each \mathbf{x}_n :

- (1) Select one of three Gaussians
probability π_k for each Gaussian
(where $\sum_k \pi_k = 1$)
- (2) Sample \mathbf{x}_n from that Gaussian



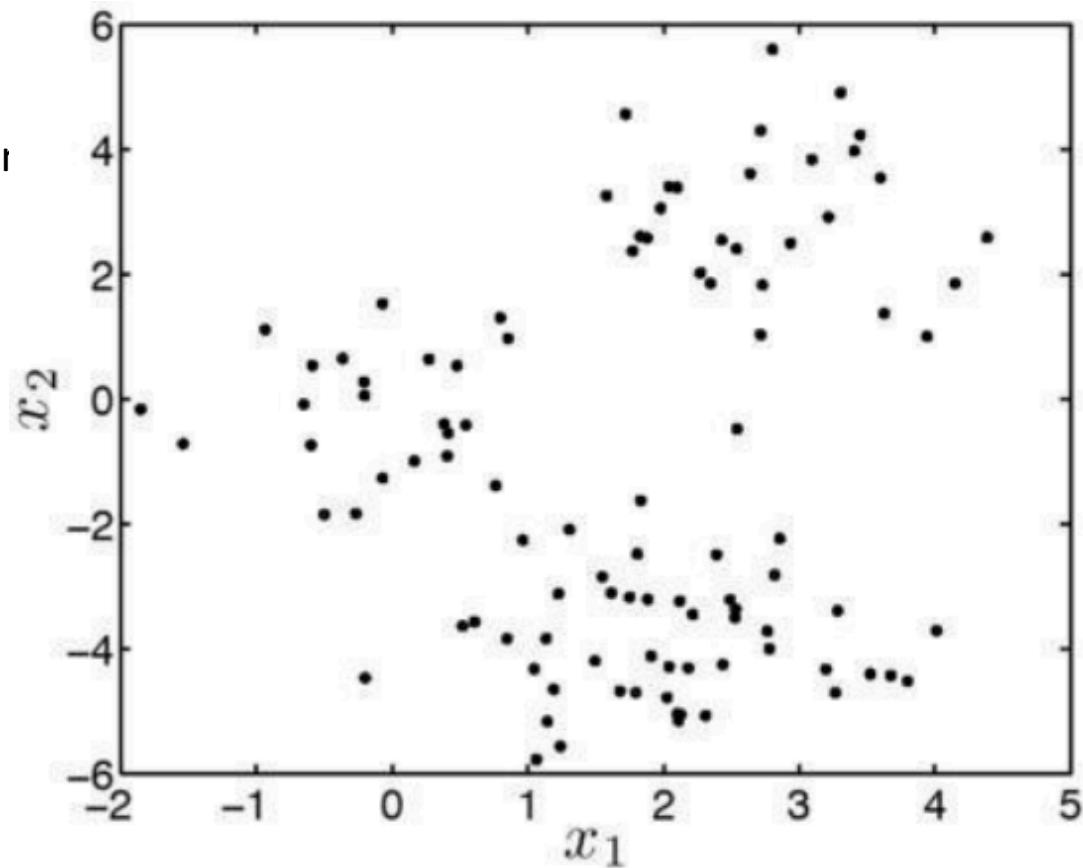
The Generative Picture (again)

- How could we *generate* this data?

For each \mathbf{x}_n :

- (1) Select one of three Gaussians probability π_k for each Gaussian (where $\sum_k \pi_k = 1$)
- (2) Sample \mathbf{x}_n from that Gaussian

- Use $z_{nk} = 1$ to mean individual n was “sampled from” generator k ($z_{nj} = 0$ for all other $j \neq k$)
- Each Gaussian is modeled with mean and covariance μ_k and Σ_k

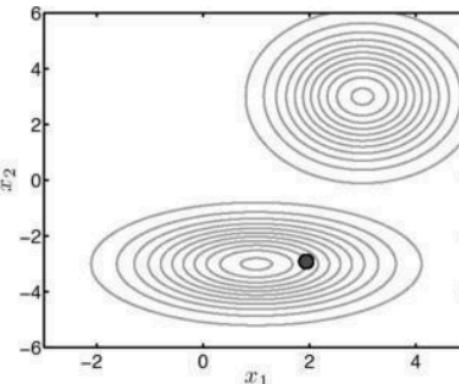


The Generative Picture (again)

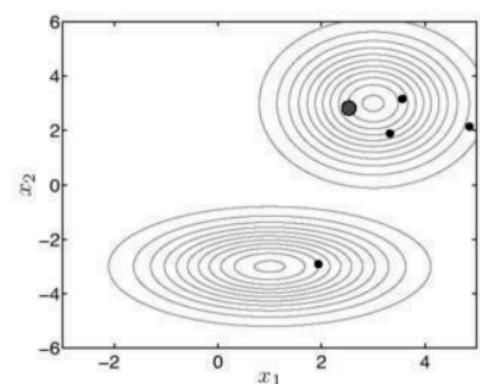
- How could we *generate* this data?

For each \mathbf{x}_n :

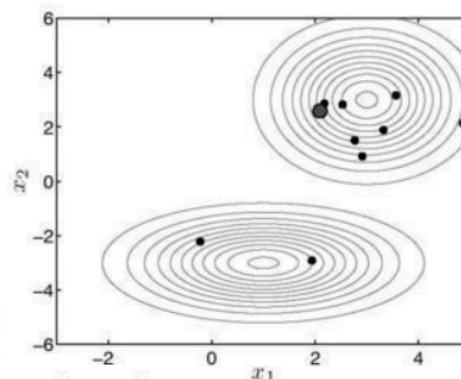
- (1) Select one of three Gaussians probability π_k for each Gaussian (where $\sum_k \pi_k = 1$)
- (2) Sample \mathbf{x}_n from that Gaussian



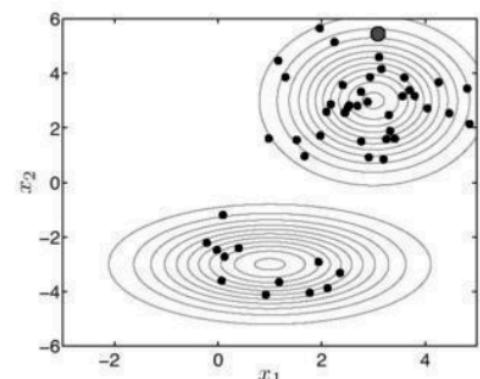
(a) The first object



(b) The first five objects



The first 10 objects



(d) The first 50 objects

$$p(\mathbf{x}_n | z_{nk} = 1, \mu_k, \Sigma_k) = \mathcal{N}(\mu_k, \Sigma_k)$$

$$\mu_1 = [3, 3]^\top, \Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\mu_2 = [1, -3]^\top, \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\pi_1 = 0.7, \pi_2 = 0.3$$

Note axis scale; x_2 is being squashed

The EM Algorithm

- Our learning task: infer, from observed data, the component
 - parameters (μ_k, Σ_k) and π_k , **and**
 - assignments z_{nk} of objects to components.
- Similar to K-means problem:
 - component **parameters** depend on **assignment**,
 - and **assignments** depend on **parameters**.
- In this probabilistic framework, we also have a two-step algorithm that alternates between steps until convergence; but now the steps are defined in terms of calculating
 - **expectations** (to update data-to-cluster **assignments**) and then
 - making adjustments to the parameters that **maximize** the likelihood (update **cluster definitions**: **parameters**)
- The **Expectation Maximization (EM)** algorithm.

Derive: the (general) Mixture Model Likelihood

Likelihood of getting individual \mathbf{x}_n assuming it was generated from cluster k : $p(\mathbf{x}_n | z_{nk} = 1, \Delta_k)$

Δ_k Represents the parameters of the k th density (In a moment, we'll use a Gaussian distribution, but could be any suitable density)

$\Delta = \{\Delta_1, \dots, \Delta_K\}$ Collection of parameters for all of the mixture components

$\pi = \{\pi_1, \dots, \pi_K\}$ Collection of all of the probabilities of the mixture components

Or goal: find the likelihood of the data object \mathbf{x}_n under the **whole model**: $p(\mathbf{x}_n | \Delta, \pi)$

Start with likelihood for one cluster:

$$p(\mathbf{x}_n | z_{nk} = 1, \Delta)$$

Need to “get rid” of z_{nk} (i.e., over all mixtures)

Multiply both sides by the probability that object n is in cluster k :

$$p(\mathbf{x}_n | z_{nk} = 1, \Delta)p(z_{nk} = 1) = p(\mathbf{x}_n | \Delta_k)p(z_{nk} = 1)$$

The probability chain rule: $P(a|b,c) \times P(b) = P(a,b|c)$

By definition: $p(z_{nk} = 1) = \pi_k$

$$p(\mathbf{x}_n, z_{nk} = 1 | \Delta, \pi) = p(\mathbf{x}_n | \Delta_k)\pi_k$$

Marginalize over all of the individual components:

$$\sum_{k=1}^K p(\mathbf{x}_n, z_{nk} = 1 | \Delta, \pi) = \sum_{k=1}^K p(\mathbf{x}_n | \Delta_k)\pi_k$$

$$p(\mathbf{x}_n | \Delta, \pi) = \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \Delta_k)$$

Make standard assumption that all data points are independent (given their mixture):

This is the likelihood of all N data points

$$p(\mathbf{X} | \Delta, \pi) = \prod_{n=1}^N \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \Delta_k)$$

$$p(\mathbf{X}|\Delta, \boldsymbol{\pi}) = \prod_{n=1}^N \sum_{k=1}^K \pi_k p(\mathbf{x}_n|\Delta_k)$$

Maximizing the Mixture Model Likelihood

- Now we'll look at an instance of the **EM** algorithm for Gaussian mixtures: a **Gaussian Mixture Model**.
- We'll want to do maximization, so easier to work with the logarithm of the likelihood

$$L = \log p(\mathbf{X}|\Delta, \boldsymbol{\pi}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

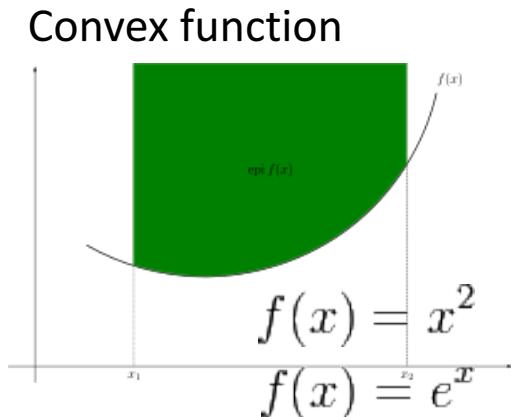
- Immediate problem!**: the summation inside the log makes finding the optimal parameters challenging

We want to take partial derivatives w.r.t. $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$

- Trick**: derive a lower-bound on L and maximize *that*

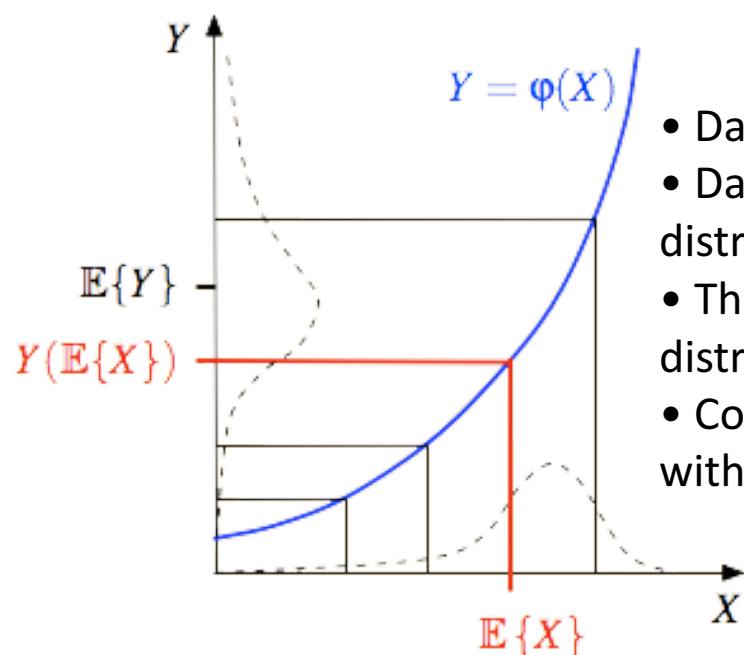
Jensen's Inequality

- A very general result that has wide application in convex optimization and probability theory
- Generally: relates the value of a convex function of an integral to the integral of the convex function
- Simplest probabilistic form: **convex transformation of a mean is less than or equal to the mean after convex transformation**
- If φ is a convex function: $\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$.



A visual proof (for probabilistic case):

(Concave functions, e.g., $\log(x)$, just reverse the inequality: $\varphi(\mathbb{E}[X]) \geq \mathbb{E}[\varphi(X)]$)



- Dashed curve along X axis is the hypothetical distribution of X .
- Dashed curve along Y axis is corresponding convex-mapped distribution of Y values ($Y = \varphi(X)$).
- The convex mapping $Y(X)$ increasingly “stretches” the distribution for increasing values of X .
- Consequently, the expectation of Y will always shift upwards with respect to the position of $\varphi(\mathbb{E}\{X\})$, thus:

$$\mathbb{E}\{Y\} = \mathbb{E}\{\varphi(X)\} \geq \varphi(\mathbb{E}\{X\}),$$

$$L = \log p(\mathbf{X}|\Delta, \boldsymbol{\pi}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad \text{Original log-likelihood}$$

$$\log \mathbf{E}_{p(z)} \{f(z)\} \geq \mathbf{E}_{p(z)} \{\log f(z)\} \quad \text{Jensen's inequality ('concave' form)}$$

- Need to make right-hand side of log-likelihood look like the log of an *expectation*.

– (Note: it sort of does now with π_k , except that we want to keep π_k around in order to maximize w.r.t. it!)

- Multiply *and* divide the expression in the summation over k by a new variable, q_{nk}

$$L = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \frac{q_{nk}}{q_{nk}}$$

- Restrict q_{nk} to be positive and sum to 1 over k

– I.e., q_{nk} is some probability distribution over the K components for the n th object (like a “soft” version of the z_{nk} membership function)

- Now rewrite the above as an expectation w.r.t. q_{nk} :

$$L = \sum_{n=1}^N \log \sum_{k=1}^K q_{nk} \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} = \sum_{n=1}^N \log \mathbf{E}_{q_{nk}} \left\{ \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right\}$$

- Now apply Jensen's inequality to get the lower bound we desire:

$$L = \sum_{n=1}^N \log \mathbf{E}_{q_{nk}} \left\{ \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right\} \geq \boxed{\sum_{n=1}^N \mathbf{E}_{q_{nk}} \left\{ \log \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right\}}$$

B

More Algebra to make it nicer...

$$\begin{aligned} \mathcal{B} &= \sum_{n=1}^N \mathbf{E}_{q_{nk}} \left\{ \log \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right\} \\ &= \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \left(\frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{q_{nk}} \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}. \end{aligned}$$

(We just used the expectation form in order to make use of Jensen's inequality...)

The parameters we now want to adjust in order to (locally) maximize B , which in turn corresponds to local maxima of L $q_{nk}, \pi, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$

Now to get down to business with the actual EM algorithm...

Where we are headed: we will find equations that are updates to the params that maximize the likelihood (via B). – so maximize B w.r.t. the parameters.

But, π_k , $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ will turn out to be in terms of q_{nk}

So make EM an iterative algorithm (analogous to K-means):

Update π_k , $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ with q_{nk} (i.e., the assignments) fixed (**M** step: maximize)

Update q_{nk} (by taking the expectation w.r.t. the unknown q_{nk} assignments – **E** step)

Maximizing B

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}$$

Find values of $q_{nk}, \pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ that correspond to a local maxima of

Update for π_k

The “mixture” probability: $\sum_k \pi_k = 1$

Thus, optimization w.r.t. π_k is constrained

Use Lagrange terms to capture constraint! Add it to B (still only involving π_k)

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k - \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) + \dots$$

Take partial derivative w.r.t. π_k and set to 0

$$\frac{\partial \mathcal{B}}{\partial \pi_k} = \frac{\sum_{n=1}^N q_{nk}}{\pi_k} - \lambda = 0$$

$$\sum_{n=1}^N q_{nk} = \lambda \pi_k$$

Now have a λ – solve for it:
Sum both sides over k , b/c that eliminates π_k and q_{nk}

$$\sum_{k=1}^K \sum_{n=1}^N q_{nk} = \lambda \sum_{k=1}^K \pi_k$$

$$\begin{aligned} \sum_{n=1}^N 1 &= \lambda \\ \lambda &= N \end{aligned}$$

Plug in N for λ – now have π_k :

$$\pi_k = \frac{1}{N} \sum_{n=1}^N q_{nk}$$

Yay! One down.

Maximizing B

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \boxed{\sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}$$

Find values of $q_{nk}, \pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ that correspond to a local maxima of

Update for $\boldsymbol{\mu}_k$

$p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ in this case will be a multivariate Gaussian, so rewrite B

$$\begin{aligned} \mathcal{B} &= \dots + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \left(\frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right) \right) + \dots \\ &= \dots - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log ((2\pi)^d |\boldsymbol{\Sigma}_k|) - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) + \dots \end{aligned}$$

Set to 0 and solve for $\boldsymbol{\mu}_k$

$$\sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0$$

$$\sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} \mathbf{x}_n = \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k$$

$$\sum_{n=1}^N q_{nk} \mathbf{x}_n = \boldsymbol{\mu}_k \sum_{n=1}^N q_{nk}$$

Two down.

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N q_{nk} \mathbf{x}_n}{\sum_{n=1}^N q_{nk}}$$

Take partial derivative w.r.t. $\boldsymbol{\mu}_k$ and set to 0

Use this linear algebra derivative fact:

$$f(\mathbf{w}) = \mathbf{w}^\top \mathbf{C} \mathbf{w}, \quad \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{C}\mathbf{w} \quad \text{... and chain rule of derivatives:}$$

$$\begin{aligned} \frac{\partial \mathcal{B}}{\partial \boldsymbol{\mu}_k} &= -\frac{1}{2} \sum_{n=1}^N q_{nk} \times \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)} \times \frac{\partial (\mathbf{x}_n - \boldsymbol{\mu}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k). \end{aligned}$$

Maximizing B

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \boxed{\sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}$$

Find values of $q_{nk}, \pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ that correspond to a local maxima of

Update for $\boldsymbol{\Sigma}_k$

... also only shows up in the term with $p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

$$\mathcal{B} = \dots -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \left((2\pi)^d |\boldsymbol{\Sigma}_k| \right) - \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) + \dots$$

Take partial derivative w.r.t. $\boldsymbol{\Sigma}_k$ and set to 0

Use these two linear algebra derivative facts: $\frac{\partial \log |\mathbf{C}|}{\partial \mathbf{C}} = (\mathbf{C}^T)^{-1}$ $\frac{\partial \mathbf{a}^T \mathbf{C}^{-1} \mathbf{b}}{\partial \mathbf{C}} = -(\mathbf{C}^T)^{-1} \mathbf{a} \mathbf{b}^T (\mathbf{C}^T)^{-1}$

$$\frac{\partial \mathcal{B}}{\partial \boldsymbol{\Sigma}_k} = -\frac{1}{2} \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} + \frac{1}{2} \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}$$

Note:
 $\boldsymbol{\Sigma}_k^T = \boldsymbol{\Sigma}_k$

Set to 0 and solve for $\boldsymbol{\Sigma}_k$

$$-\frac{1}{2} \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} + \frac{1}{2} \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} = 0$$

$$\frac{1}{2} \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} = \frac{1}{2} \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}$$

Pre and post multiply by $\boldsymbol{\Sigma}_k$ removes all $\boldsymbol{\Sigma}_k^{-1}$

$$\boldsymbol{\Sigma}_k \sum_{n=1}^N q_{nk} \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}_k \boldsymbol{\Sigma}_k^{-1} \sum_{n=1}^N q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\Sigma}_k$$

$$\boldsymbol{\Sigma}_k \sum_{n=1}^N q_{nk} = \sum_{n=1}^N q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N q_{nk}}$$

Three down.

Maximizing B

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}$$

Find values of $q_{nk}, \pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ that correspond to a local maxima of

Update for q_{nk}

Shows up in all three terms! And has this constraint: $\sum_{k=1}^K q_{nk} = 1$ Lagrangian term for constraint

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk} - \lambda \left(\sum_{k=1}^K q_{nk} - 1 \right)$$

Take partial derivative w.r.t. q_{nk} and set to 0

Need the derivative product rule: for $f(a) = g(a)h(a)$

$$\frac{\partial f(a)}{\partial a} = g(a) \frac{\partial h(a)}{\partial a} + \frac{\partial g(a)}{\partial a} h(a)$$

$$\frac{\partial \mathcal{B}}{\partial q_{nk}} = \log \pi_k + \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - (1 + \log q_{nk}) - \lambda$$

Set to 0, rearrange, exponentiate, and solve for q_{nk}

$$1 + \log q_{nk} + \lambda = \log \pi_k + \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\exp(\log q_{nk} + (\lambda + 1)) = \exp(\log \pi_k + \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))$$

$$q_{nk} \exp(\lambda + 1) = \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Same trick we used for $\boldsymbol{\mu}_k$: sum over k on both sides makes q_{nk} go to 1 on the left side:

$$\exp(\lambda + 1) \sum_{k=1}^K q_{nk} = \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\exp(\lambda + 1) = \sum_{k=1}^K \pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

Can substitute this here and solve for q_{nk} :

$$q_{nk} = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Done!

Expectation Maximization (EM) for GMM

$$\mathcal{B} = \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) - \sum_{n=1}^N \sum_{k=1}^K q_{nk} \log q_{nk}$$

The final expressions:

$$\pi_k = \frac{1}{N} \sum_{n=1}^N q_{nk}$$

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N q_{nk} \mathbf{x}_n}{\sum_{n=1}^N q_{nk}}$$

$$\boldsymbol{\mu}_k = \frac{\sum_n z_{nk} \mathbf{x}_n}{\sum_n z_{nk}}$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N q_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N q_{nk}}$$

$$q_{nk} = \frac{\pi_k p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Posterior probability of object n belonging to cluster k

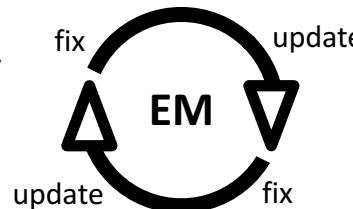
$$p(z_{nk} = 1 | \mathbf{x}_n, \boldsymbol{\pi}, \boldsymbol{\Delta}) = \frac{p(z_{nk} = 1 | \pi_k) p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K p(z_{nj} = 1 | \pi_j) p(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} = q_{nk}.$$

Cluster membership probability

$$q_{nk}$$

E step

(expected value of unknown assignments z_{nk})



All about the model

$\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \boldsymbol{\pi}$
M step

(maximizing w.r.t cluster membership)

The mean value of q_{nk} for a particular cluster k :

Average of all posterior probabilities of belonging to cluster k
i.e., the expected proportion of the data belonging to cluster k

Average of the data objects weighted by q_{nk} :

When all cluster membership is such that posterior prob's are 0 or 1,
then this is just the proportion of the data assigned to component k
weighted covariance :

Looks like Bayes' Rule! :

Running EM with a GMM

- Initialize the parameters:
 - Mixture parameters: random means and covariances
 - Mixture priors: could choose uniform $\pi_k = 1/K$

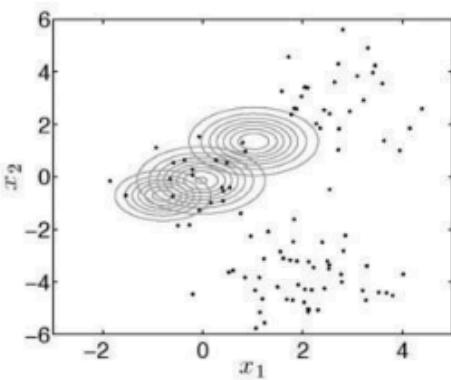
To find cluster membership:

q_{nk} = posterior prob of n belonging to k

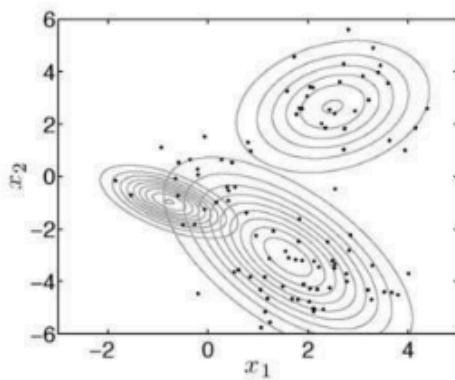
Hard assignment: for each indiv. n , choose the k with the highest q_{nk}

Or, consider the distribution – reveals interesting relationships of individual to more than one cluster; e.g.,

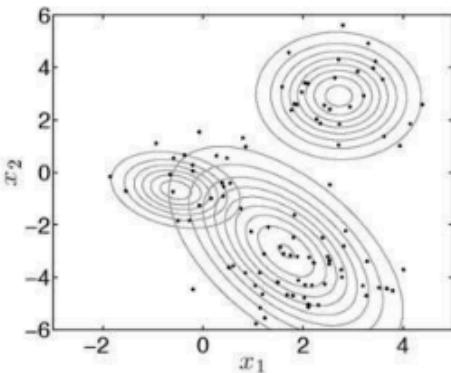
$$q_{n1} = 0.53, \quad q_{n2} = 0.45, \quad q_{n3} = 0.02$$



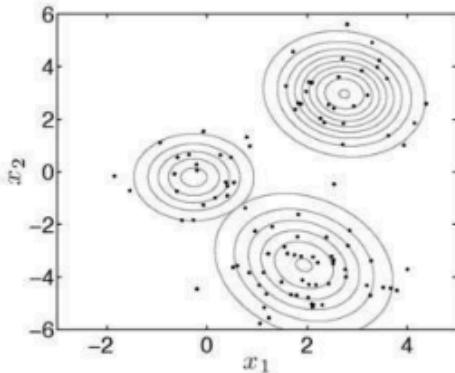
(a) The three randomly initialised Gaussian mixture components



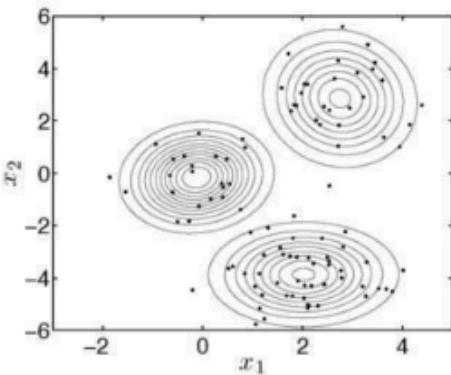
(b) The three components after one iteration of the EM algorithm



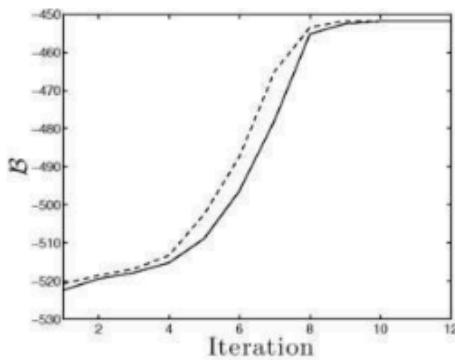
(c) The three components after five iterations of the EM algorithm



(d) The three components after seven iterations of the EM algorithm

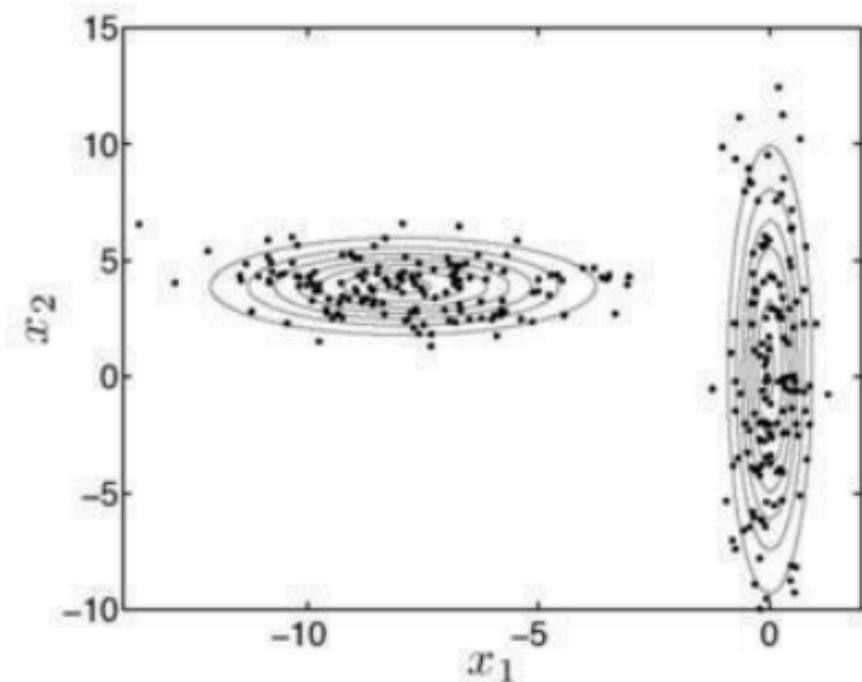
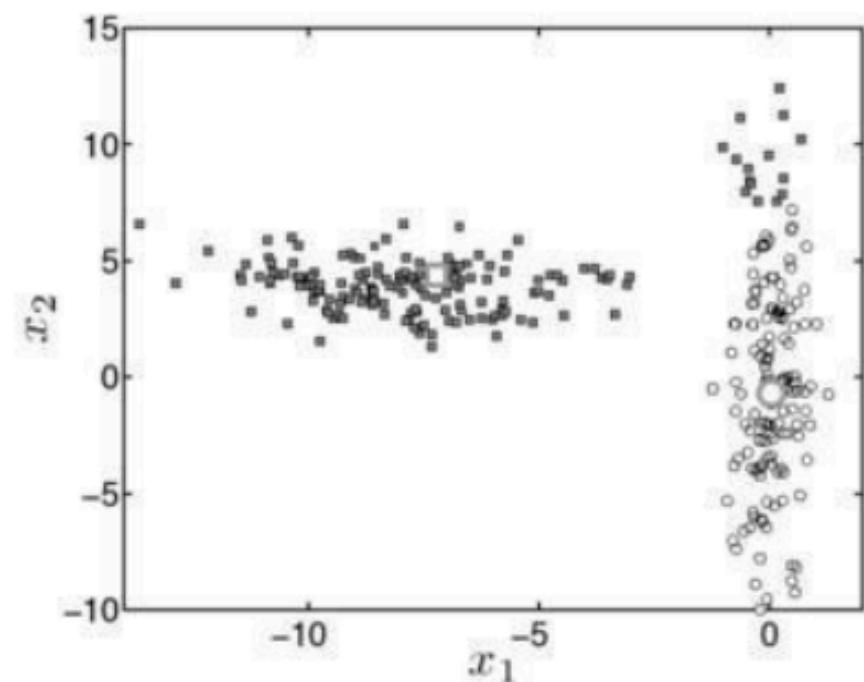


(e) The three components at convergence of the EM algorithm



(f) The evolution of the bound B (solid line, Equation 6.8) and log likelihood L (dashed line, Equation 6.5)

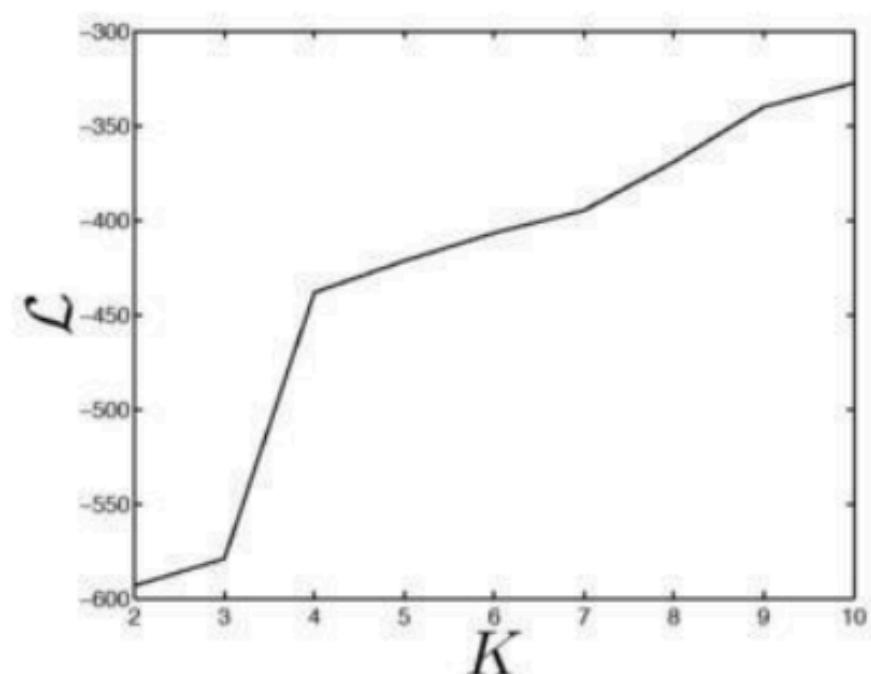
GMM does well on problem hard for K-means



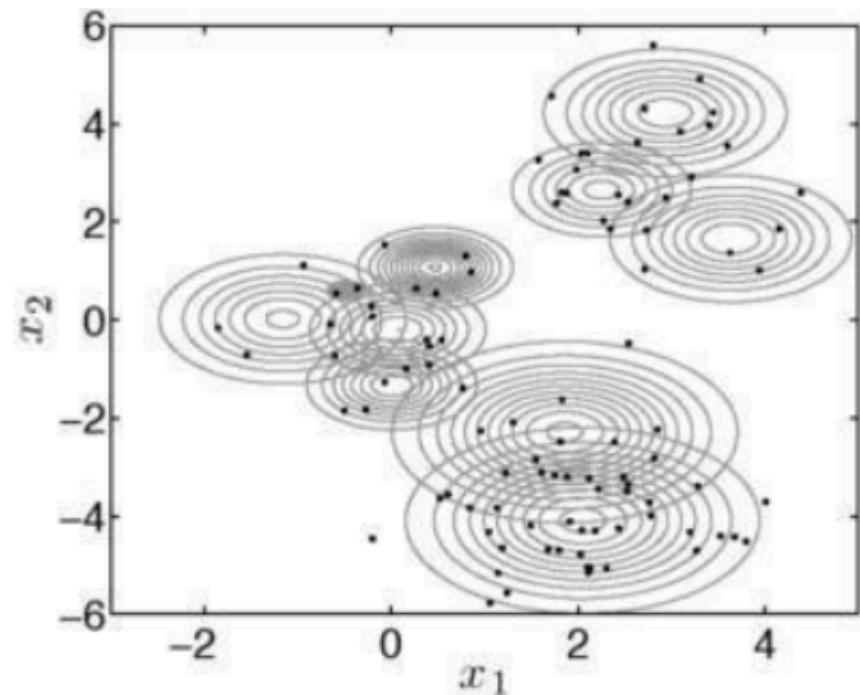
GMM with K=2

Choosing K

- Similar to K -means, where we can't just choose the K that minimizes the total distances ($K=N \rightarrow D=0$), we can't just choose the K that maximizes the log likelihood L (or bound B); it increases with more mixtures:



(a) The increase in model likelihood as the number of components increases



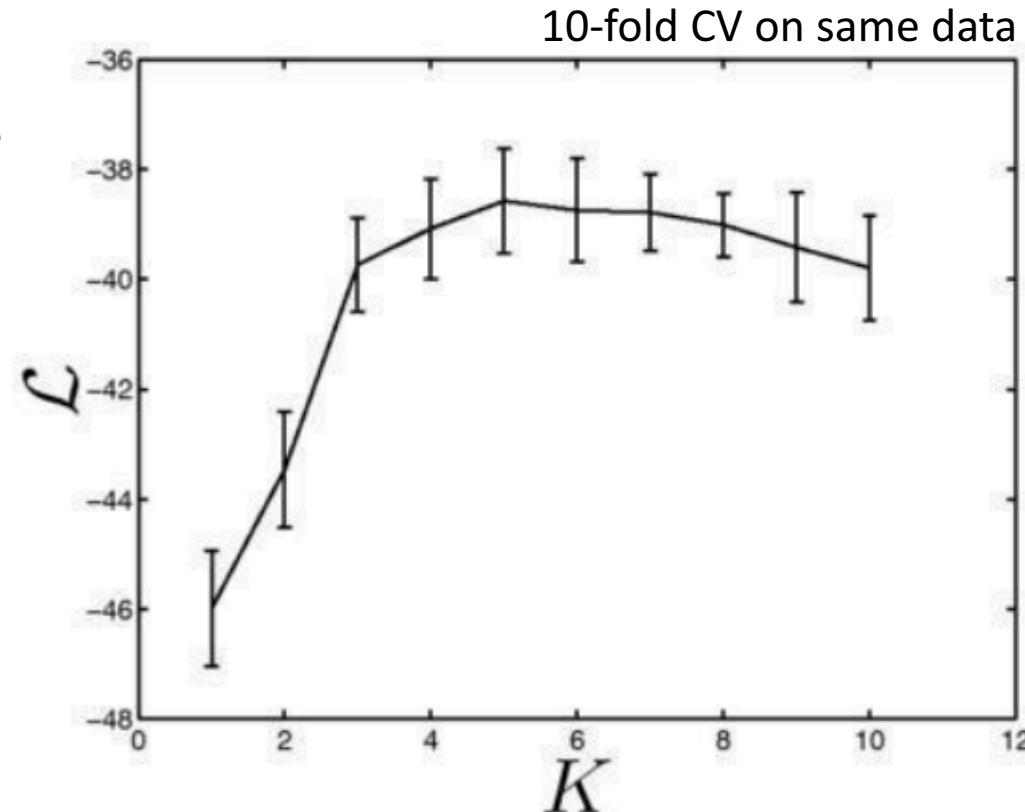
(b) An example of the model at convergence for $K = 10$

The power of a generative model

- **However:** because mixture models are generative models, we can run cross-validation:
 - For each potential K : Hold out data, fit mixtures, then measure likelihood of held-out data

It's not perfect (somewhere between 3 and 8)

But with non-generative K-means, we have no direct basis for choosing K



Other Mixtures besides Gaussians

- Could be any probability density

$$p(\mathbf{X}|\Delta, \pi) = \prod_{n=1}^N \sum_{k=1}^K \pi_k p(\mathbf{x}_n|\Delta_k)$$

- For D=10 dimensional binary data, e.g.,

$$\mathbf{x}_n = [0, 1, 0, 1, 1, 1, 0, 0, 0, 1]$$

- Represent as product of (i.e., independent) Bernoulli distributions:

$$p(\mathbf{x}_n|\mathbf{p}_k) = \prod_{d=1}^D p_{kd}^{x_{nd}} (1 - p_{kd})^{1-x_{nd}}$$

Mixture probabilities

$$\mathbf{p}_k = [p_{k1}, \dots, p_{kD}]^\top$$