

Mixture Models and EM: Model-Based Clustering

CSE 446: Machine Learning

Slides by Emily Fox

University of Washington

May 22, 2019

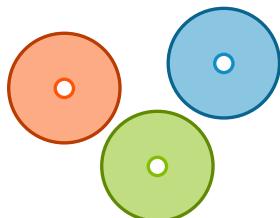
Limitations of k-means

Assigns observations to closest cluster center

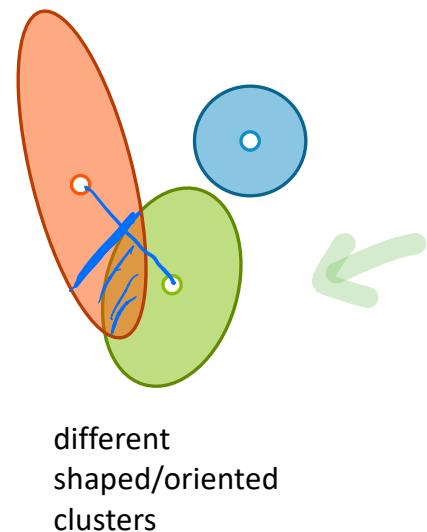
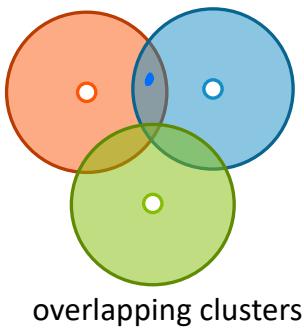
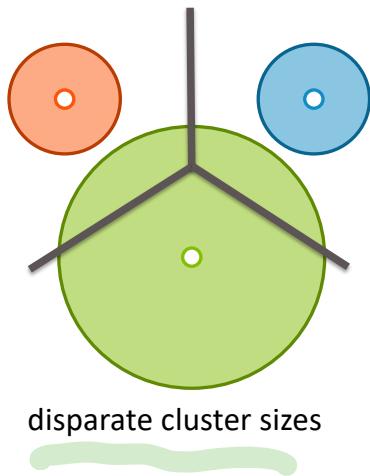
$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

Only center matters
Not cluster shapes

Equivalent to assuming
spherically symmetric clusters

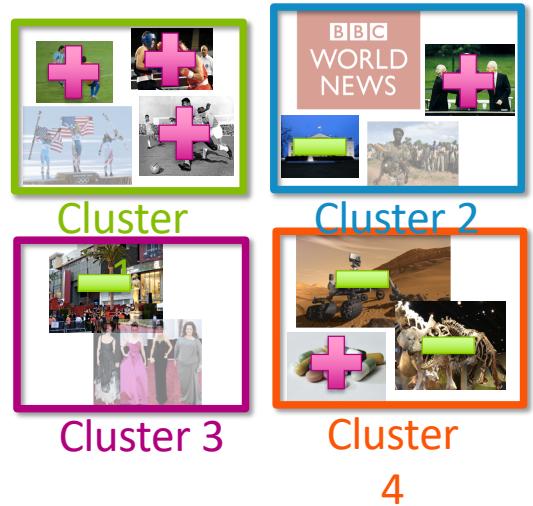


Failure modes of k-means



Motivates probabilistic model: Mixture model

- Take uncertainty in assignment into account
e.g., when clustering documents, might want to say 54% chance document is **world news**, 45% **science**, 1% **sports**, and 0% **entertainment**
- Accounts for cluster **shapes** not just **centers**

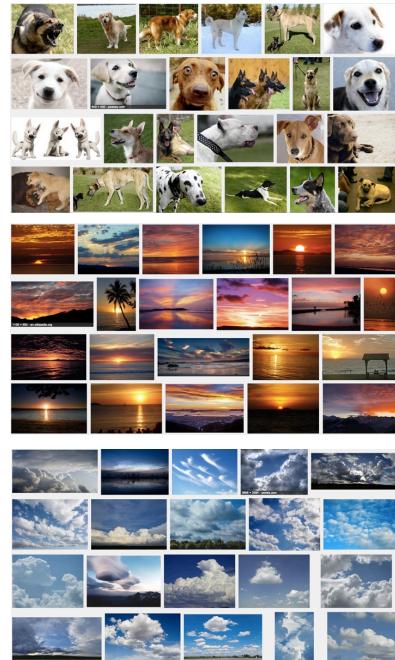
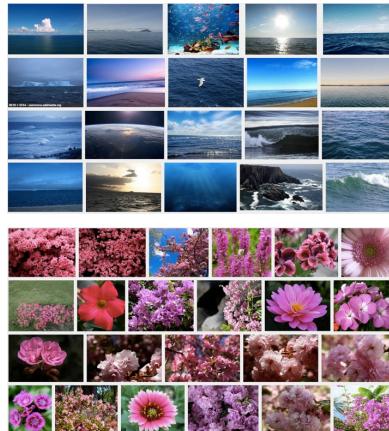


Mixture models

Motivating application: Clustering images

Discover groups of similar images

- Ocean
- Pink flower
- Dog
- Sunset
- Clouds
- ...



Simple image representation

Consider average red, green, blue pixel intensities



[R = 0.05, G = 0.7, B = 0.9]



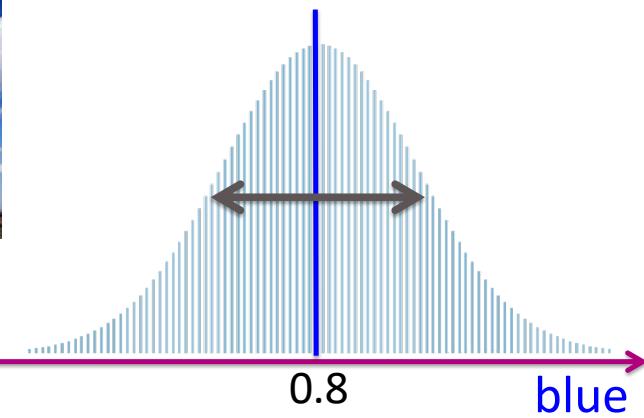
[R = 0.85, G = 0.05, B = 0.35]



[R = 0.02, G = 0.95, B = 0.4]

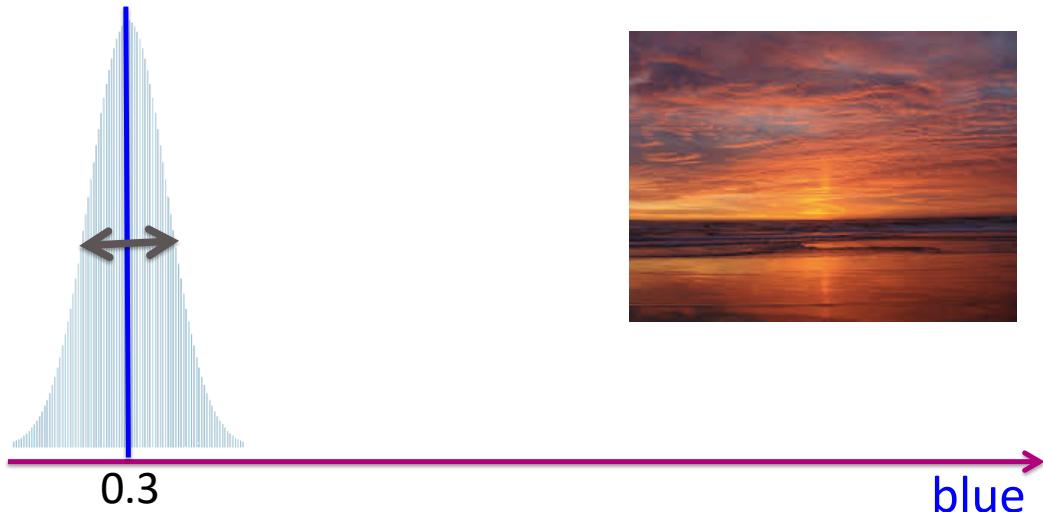
Distribution over all **cloud** images

Let's look at just the **blue** dimension



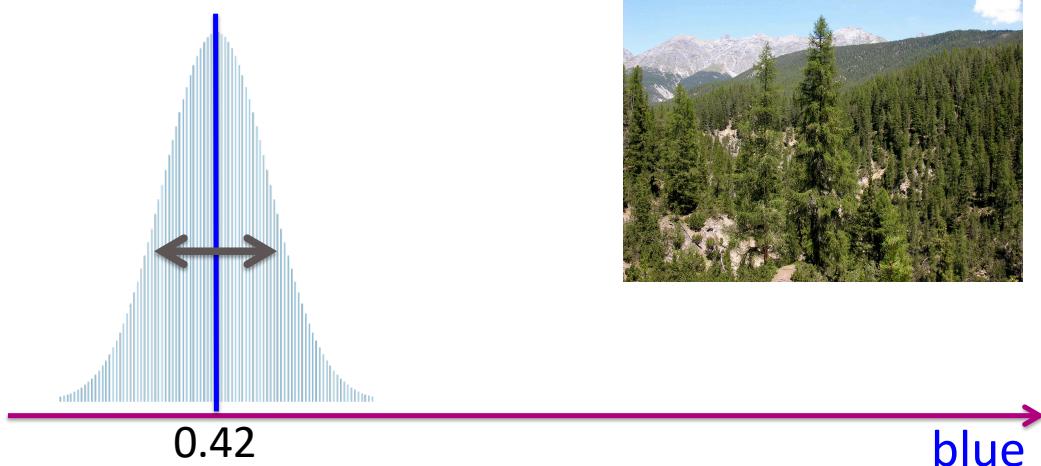
Distribution over all sunset images

Let's look at just the blue dimension

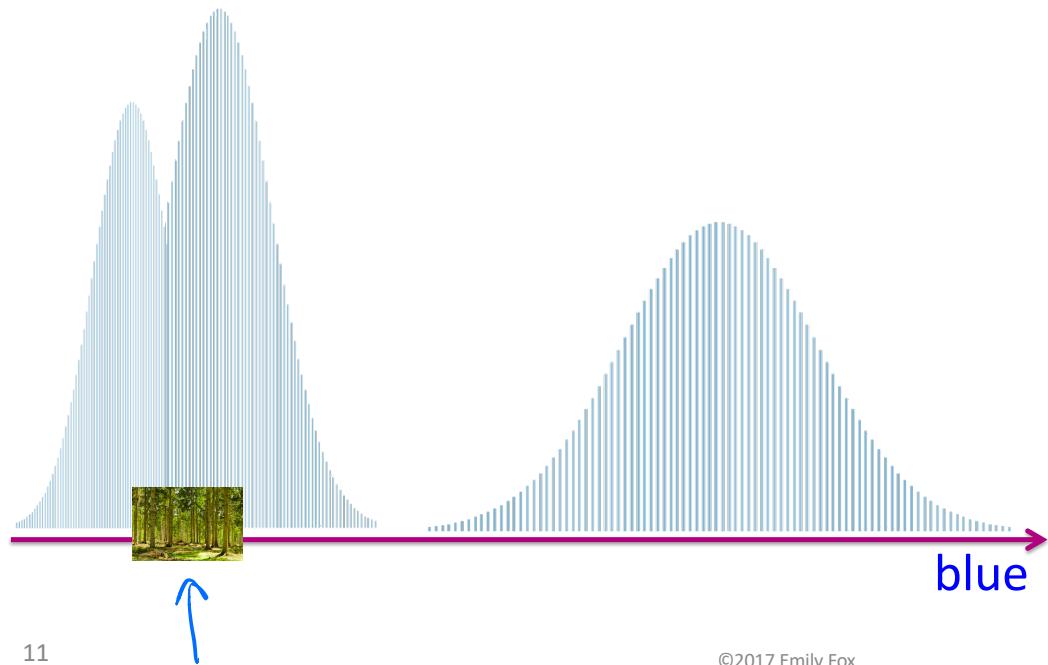


Distribution over all forest images

Let's look at just the blue dimension



Distribution over **all** images



$$\vec{x}_i \in \mathbb{R}^d$$
$$\approx$$

Can be distinguished along other dim

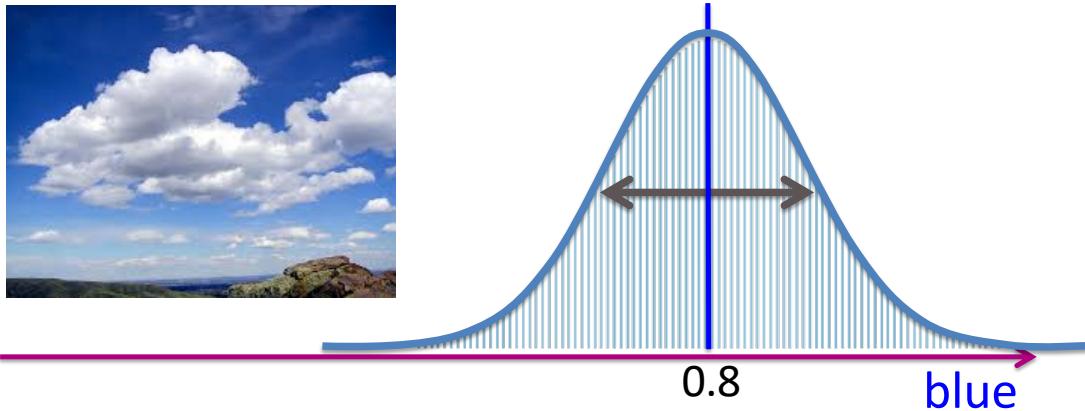
Now look at the **red** dimension



Background: Gaussian distributions

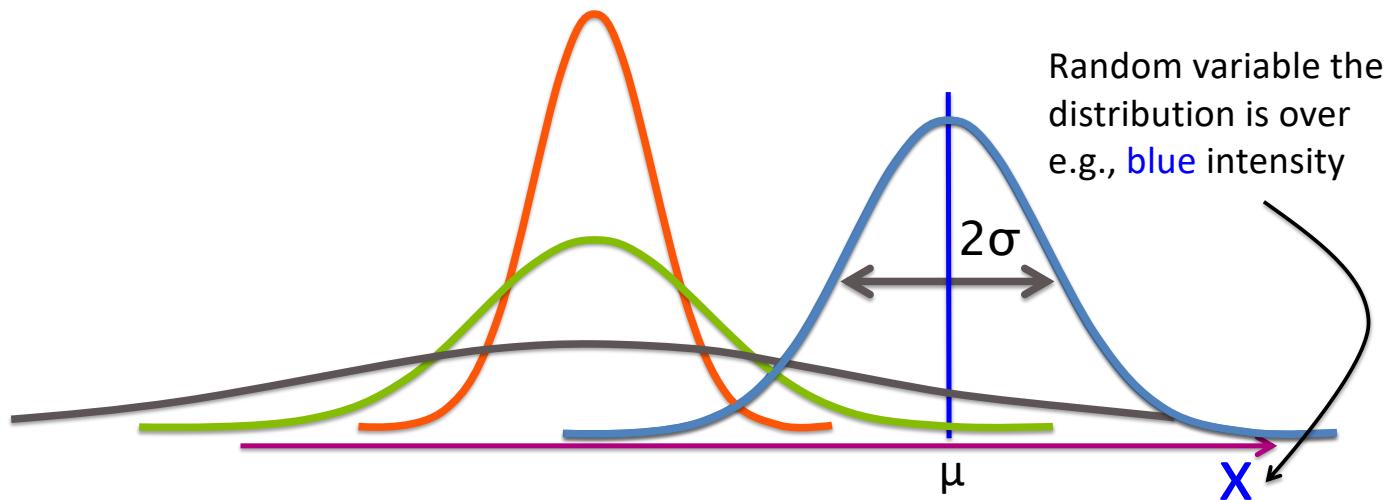
Model for a given image type

For each dim of the [R, G, B] vector, and each image type, assume a **Gaussian distribution** over color intensity



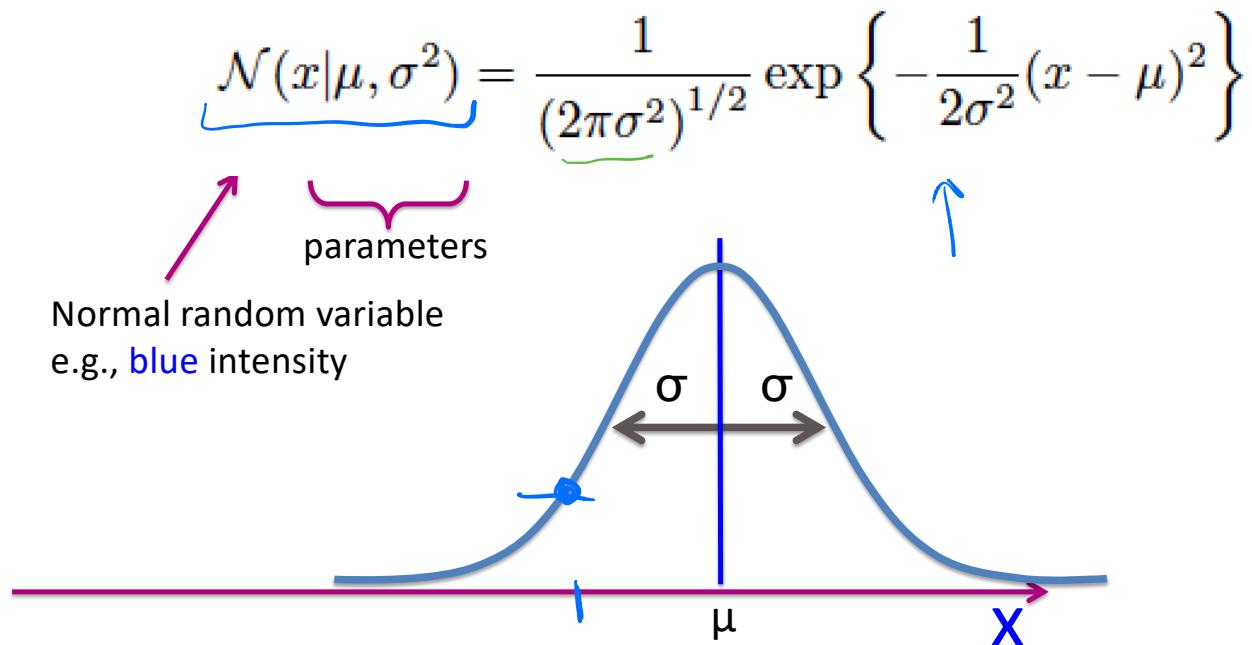
1D Gaussians

Fully specified by mean μ and variance σ^2 (or st. dev. σ)



x is sample from $N(\mu, \sigma^2)$

Notation a 1D Gaussian distribution



$$\gamma_i > 0$$

Multivariate Gaussian density

$$\Sigma_{ij} = \text{cov}(x_i, x_j)$$

$d \times d$

$$\vec{x} = (x_1, \dots, x_d)$$

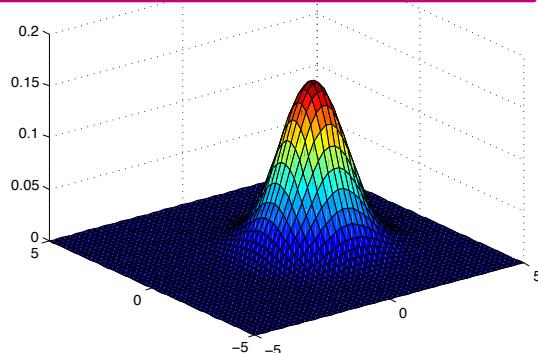
$$\vec{\mu} = (\mu_1, \dots, \mu_d)$$

$$\Sigma \in \mathbb{R}^{d \times d}$$

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$\mathcal{N}(x | \mu, \Sigma)$
 Random vector
 e.g., [R, G, B] intensities

parameters



$$|\Lambda| = \det(\Lambda)$$

$$\begin{aligned} \det(2\pi\Sigma) &= (2\pi)^d \det(\Sigma) \\ \sqrt{2\pi\Sigma} &= \sqrt{2\pi} \sqrt{\det(\Sigma)} \end{aligned}$$

A
 $x^T A x$ quadratic form

$$(x - \mu)^T \Sigma^{-1} (x - \mu)$$

Multivariate Gaussian

Fully specified by mean μ and covariance Σ

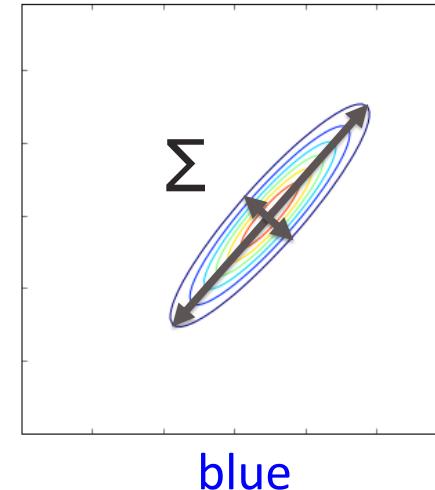
$$\mathbf{x} = (x_1, x_d)$$

$$\Sigma_{ij} = \text{cov}(x_i, x_j)$$

$$\mu = [\mu_{\text{blue}}, \mu_{\text{green}}]$$

$$\Sigma = \begin{pmatrix} \sigma_{\text{blue}}^2 & \sigma_{\text{blue},\text{green}} \\ \sigma_{\text{green},\text{blue}} & \sigma_{\text{green}}^2 \end{pmatrix}$$

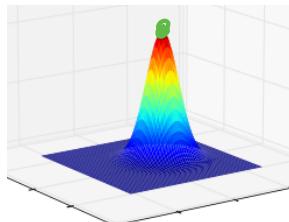
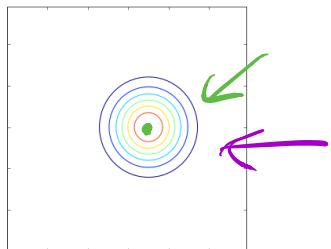
covariance determines orientation + spread



$d=2$

Covariance structure

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



for contour plot
all points where

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

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu} = 0, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} \right\}$$

quadratic

$$\boxed{\mathbf{x}^T \Sigma^{-1} \mathbf{x} = 1}$$

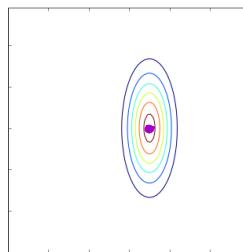
$$(\mathbf{x}_1, \mathbf{x}_2) \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{\sigma^2} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \frac{\mathbf{x}_1^2}{\sigma^2} + \frac{\mathbf{x}_2^2}{\sigma^2} = 1$$

$$\underline{\mathbf{x}_1^2 + \mathbf{x}_2^2 = \sigma^2}$$

$\mathbf{x} = (\mathbf{x}_1 \mathbf{x}_2 \dots)$ then x_i 's are indep $N(\mu_i, \sigma_i^2)$ if Σ diagonal.

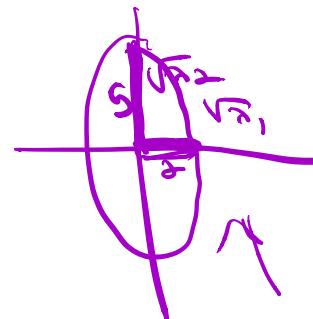
Covariance structure

$$\Sigma = \begin{pmatrix} \sigma_B^2 & 0 & 0 \\ 0 & \sigma_G^2 & 0 \\ 0 & 0 & \sigma_2^2 \end{pmatrix}$$

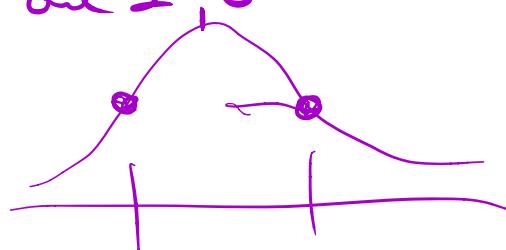
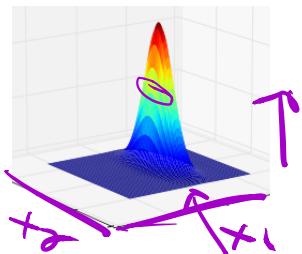


$$(x_1, x_2) \begin{pmatrix} \frac{1}{\sigma_B^2} & 0 \\ 0 & \frac{1}{\sigma_G^2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 1$$

$$\frac{x_1^2}{\sigma_B^2} + \frac{x_2^2}{\sigma_G^2} = 1$$

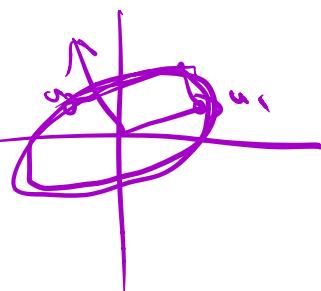
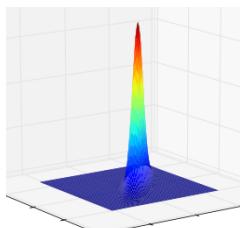
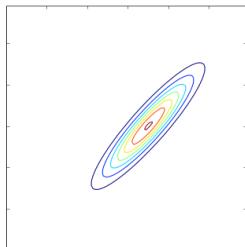


$$x_1, x_2 \text{ are } \sigma_B^2, \sigma_G^2$$



Covariance Structure

$$\Sigma = \begin{pmatrix} \sigma_B^2 & \sigma_{B,G} \\ \sigma_{G,B} & \sigma_G^2 \end{pmatrix}$$



$$x^T \Sigma^{-1} x$$

\downarrow

$$x^T U D^{-1} U^T x$$

$$y^T \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{pmatrix} y = \frac{y_1^2}{\sigma_1^2} + \frac{y_2^2}{\sigma_2^2} = 1$$

$$\Sigma = \begin{pmatrix} \sigma_1 & \sigma_{12} \\ \sigma_{21} & \sigma_2 \end{pmatrix} = U D U^T$$

$$\Sigma^{-1} = U D^{-1} U^T$$

$$UDU^T (UDU^T)^T = UDU^T (U^T D^{-1} U) = UDU^T D^{-1} U^T = U D^{-1} U^T = I$$

Important facts

- **Affine Property** $\mathbf{X} \sim N(\boldsymbol{\mu}, \Sigma)$
 $A\mathbf{X} + \mathbf{b} \sim N(A\boldsymbol{\mu} + \mathbf{b}, A\Sigma A^T)$

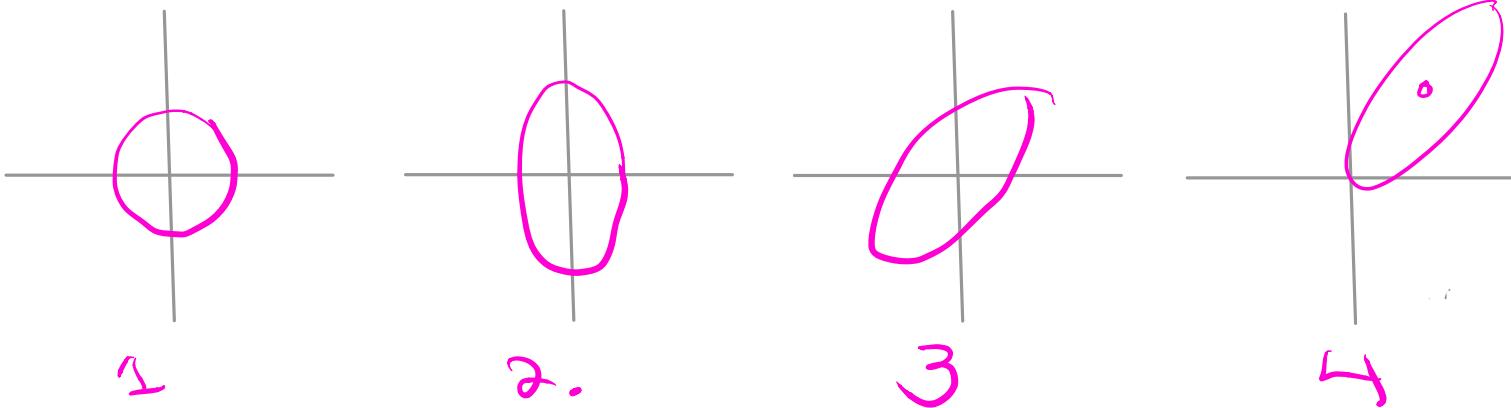
- **Constructing:** $X_1, \dots, X_d \sim N(0, 1)$ independent. Then $\mathbf{X} \sim N(0, I)$. Then
 $A\mathbf{X} + \boldsymbol{\mu} \sim N(\boldsymbol{\mu}, \Sigma)$
 $\Sigma = AA^T$ 
- **Spherizing:** If Σ is psd, symmetric, then

$$\mathbf{Y} \sim N(\boldsymbol{\mu}, \Sigma) \quad \rightarrow \quad A^{-1}(\mathbf{Y} - \boldsymbol{\mu}) \sim N(0, I)$$

$$\Sigma = AA^T$$

More intuition

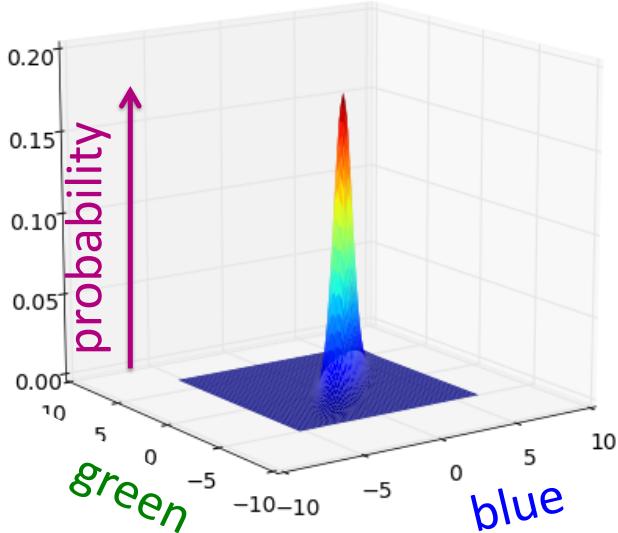
1. Start with $\underline{\mathbf{X} \sim N(0, I)}$
2. (Scaling step) $D^{1/2}\mathbf{X} \sim N(0, D)$.
3. (Rotation) $UD^{1/2}\mathbf{X} \sim N(0, \Sigma)$ where $\Sigma = UDU^T$.
4. (Translation) $UD^{1/2}\mathbf{X} + \mu \sim N(\mu, \Sigma)$



Multivariate Gaussians

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

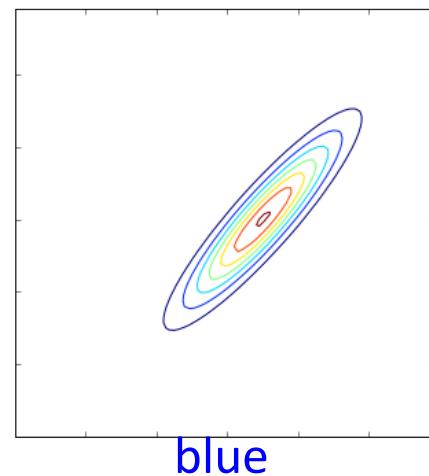
3D mesh plot



24

©2017 Emily Fox

Contour plot



CSE 446: Machine Learning

Summary

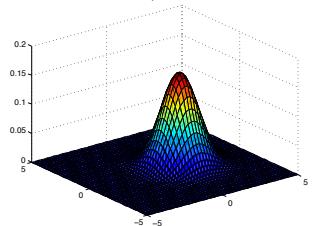
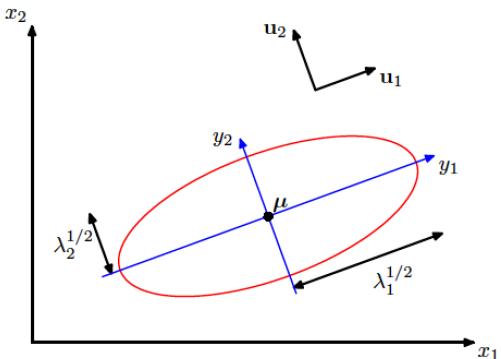
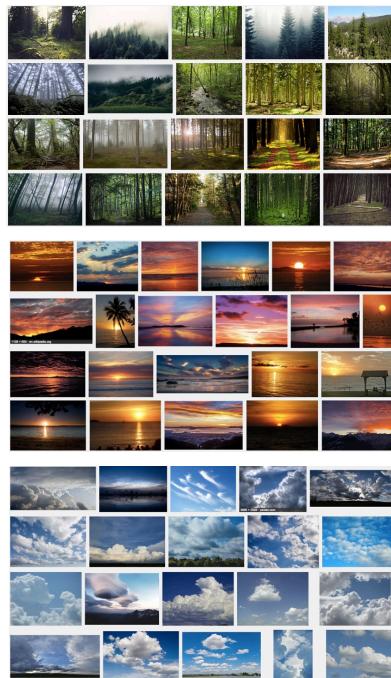
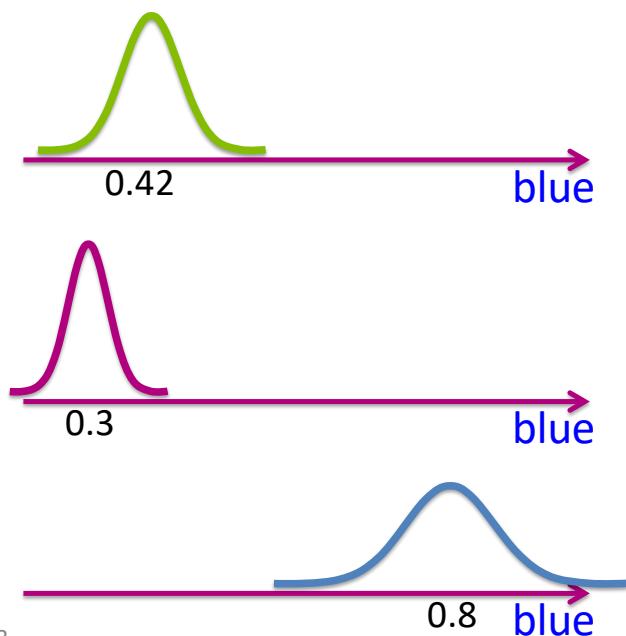


Figure 2.7 The red curve shows the elliptical surface of constant probability density for a Gaussian in a two-dimensional space $x = (x_1, x_2)$ on which the density is $\exp(-1/2)$ of its value at $x = \mu$. The major axes of the ellipse are defined by the eigenvectors u_i of the covariance matrix, with corresponding eigenvalues λ_i .



Mixture Model (to be used for clustering)

Model as Gaussian per category/cluster



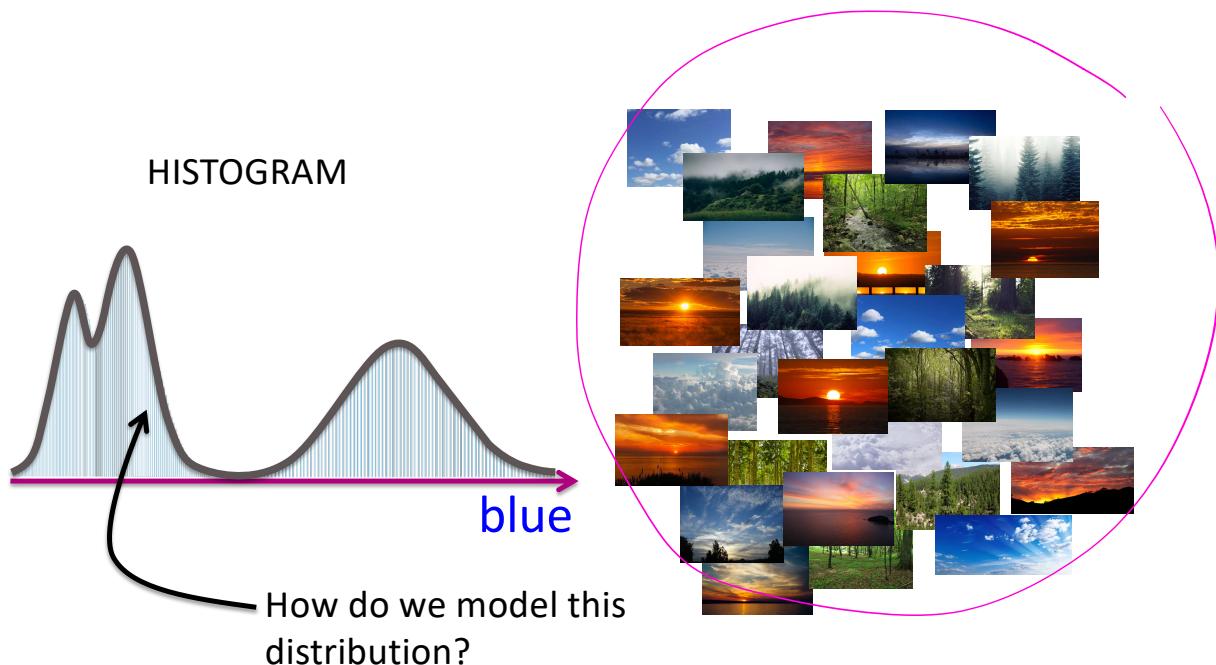
Forests

Sunsets

Clouds

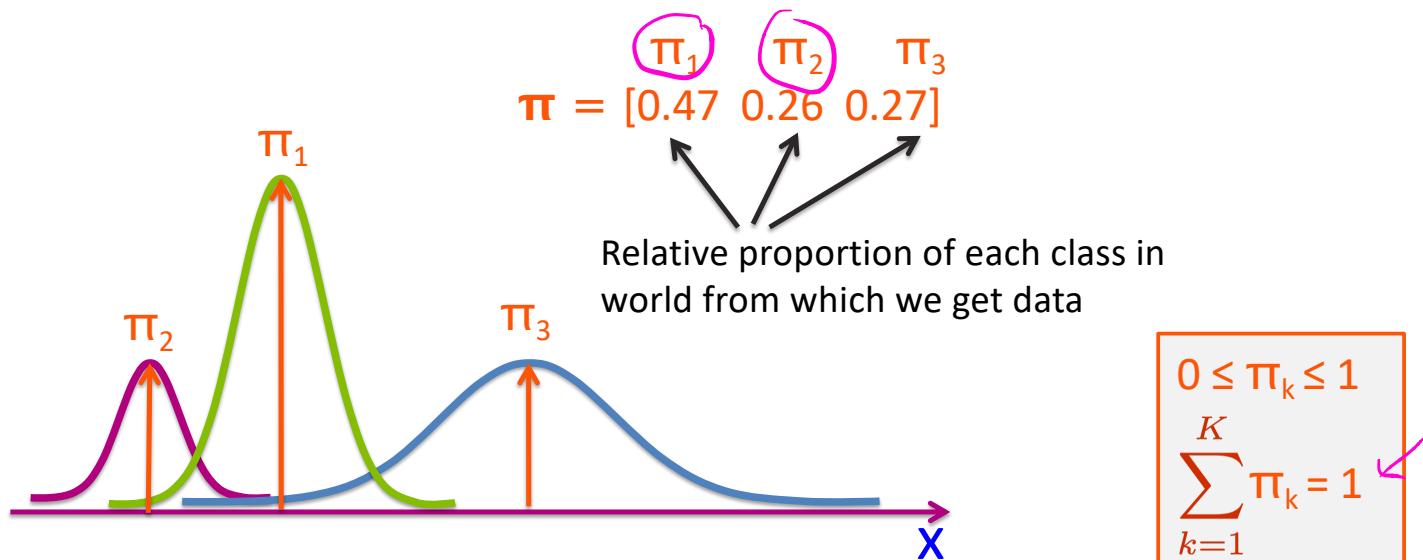
©2017 Emily Fox

Jumble of unlabeled images



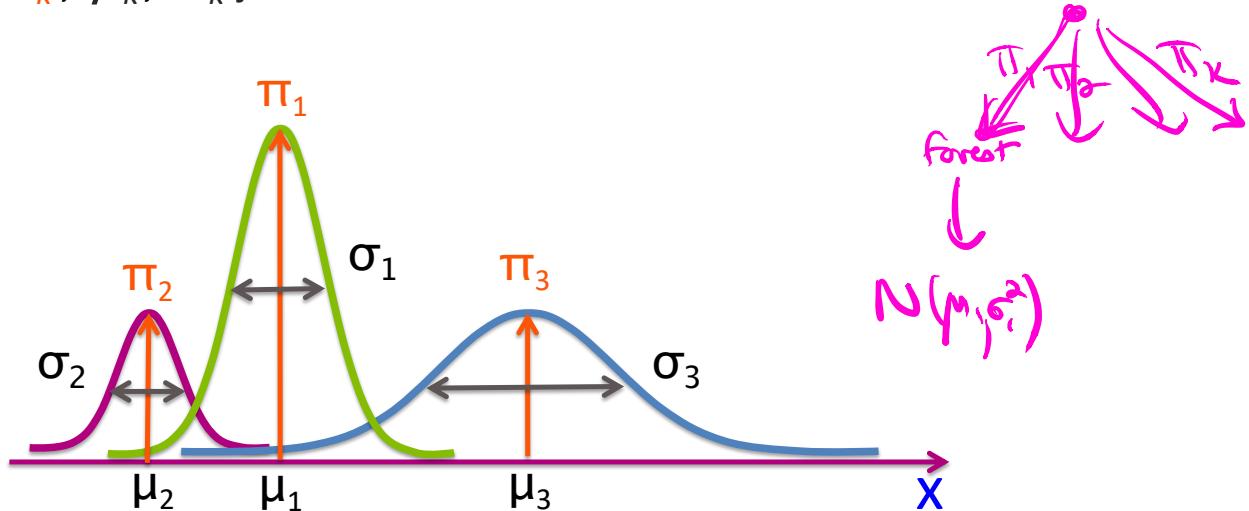
Combination of weighted Gaussians

Associate a weight π_k with each Gaussian component

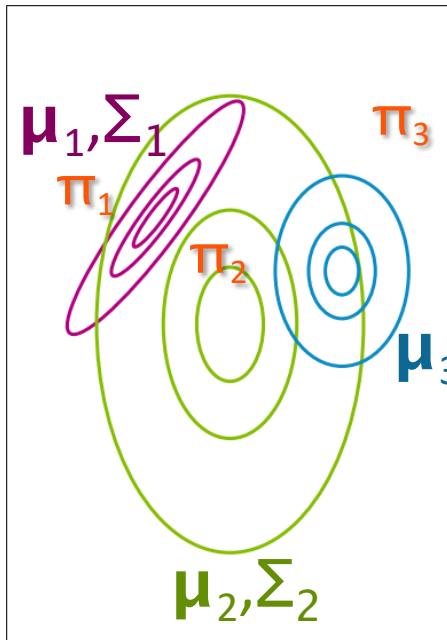


Mixture of Gaussians (1D)

Each mixture component represents a unique cluster specified by: $\{\pi_k, \mu_k, \sigma_k\}$



Mixture of Gaussians (general)



Each mixture component represents a unique cluster specified by:

$$\{\pi_k, \mu_k, \Sigma_k\}$$

↑ ↑

Mixture model

- K clusters, defined by the following parameters

$$\Theta = \{\pi_j, \mu_j, \Sigma_j\}_{j=1}^K$$

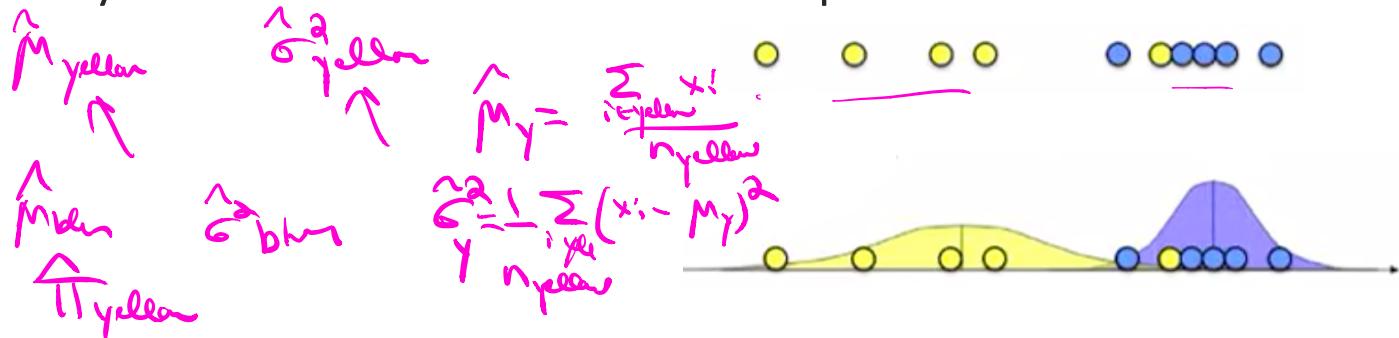
$$\sum_{j=1}^K \pi_j = 1.$$

- Problem: Assume that the data comes from such a distribution, and recover the parameters of the distribution.
- Determine, for each point, the likelihood of it belonging to cluster j, for each j.



K=2 1-D Gaussians, with unknown mean and variance

- Easy if know the source of each data point.



- What if we don't know the source?



To understand better, introduce additional “latent” variables

- K clusters, defined by the following parameters

$$\Theta = \{\pi_j, \mu_j, \Sigma_j\}_{j=1}^k \quad \sum_{j=1}^k \pi_j = 1.$$



$z_k = \begin{cases} 1 & \times \text{ custom cluster } k \\ 0 & \text{o.w} \end{cases}$

- For each point \mathbf{x} , let $\mathbf{z} = (z_1, \dots, z_K)$ indicate which cluster it was chosen from. These are called “latent variables”.

$$p(z_k = 1) = \pi_k$$

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

$$x_i \quad z_i = (z_{i1} \dots z_{iK})$$

With z 's in hand, we can compute many relevant quantities *given params*

- Conditional distribution of x given z

$$\Pr(x_i | z_{iK}=1) = N(x_i | \mu_K, \Sigma_K)$$

- Therefore:

$$\Pr(x_i) = \sum_{j=1}^K \Pr(z_{ij}=1) \Pr(x_i | z_{ij}=1) = \sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)$$

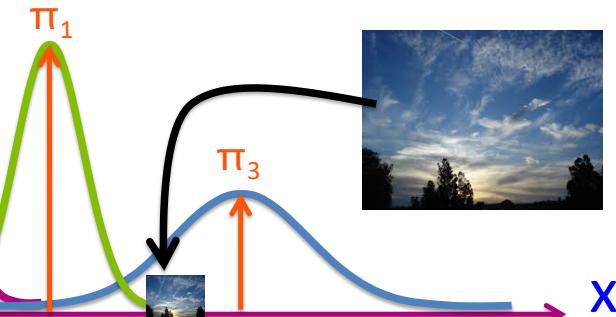
- Conditional probability

of z given x

$$\Pr(z_{iK}=1 | x_i) = \frac{\Pr(x_i | z_{iK}=1) \Pr(z_{iK}=1)}{\sum_{j=1}^K \Pr(x_i | z_{ij}=1) \Pr(z_{ij}=1)}$$

posterior

$$= \frac{\pi_K N(x_i | \mu_K, \Sigma_K)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)}$$



© 2017 Emily McKeine Learning

π_j = Prob pt coming from cluster j

Mixture model cont.

- Conditional distribution of \mathbf{x}

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Therefore:

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- We will also be interested

in the conditional probability $p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)}$
of \mathbf{z} given \mathbf{x}

$$= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}.$$



And we can now try to calculate MLE

- Given dataset from a mixture model, find parameters

$$\Theta = \{\pi_j, \mu_j, \Sigma_j\}_{j=1}^k$$

that maximize loglikelihood.

$$p(\mathbf{x}|\Theta) = \prod_{i=1}^N p(\mathbf{x}_i|\Theta)$$
$$= \sum_{k=1}^K \pi_k N(\mathbf{x}_i | \mu_k, \Sigma_k)$$

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iG})$$

MLE

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{i=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Gradient = 0 gives the following conditions:

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \mathbf{x}_i$$

$$N_k = \sum_{i=1}^N r_{ik} = \text{exp} \# \text{ptz in cluster } k$$

$$\hat{\boldsymbol{\Sigma}}_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

$$\hat{\pi}_k = \frac{N_k}{N}$$

$r_{ik} = \text{responsibility of } k^{\text{th}} \text{ cluster for point } i$
 $= \Pr(z_{ik}=1 | \mathbf{x}_i, \Theta)$

$$\begin{aligned} r_{ik} &= \underline{p(z_{ik}=1 | \mathbf{x}_i)} = \frac{p(z_k=1)p(\mathbf{x}|z_k=1)}{\sum_{j=1}^K p(z_j=1)p(\mathbf{x}|z_j=1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}. \end{aligned}$$

Does not give us a closed form $\frown\smile$

(r_{ik}) $\Pr(z_{ik}=1 | \mathbf{x}_i, \Theta)$
 \mathbf{x}_i comes from cluster k

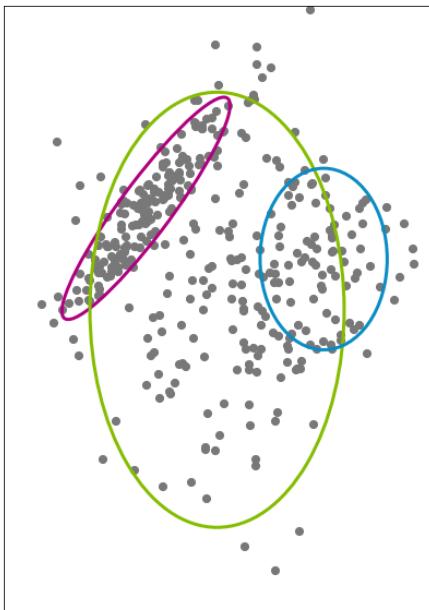
Expectation Maximization Algorithm

Two step approach based on following observation

- If we knew the z_i 's, we could estimate all the parameters.
- If we knew all the parameters we could estimate the z_i 's (or more precisely, the chance each point came from each cluster)
- EM is an iterative algorithm that alternates between these two steps.

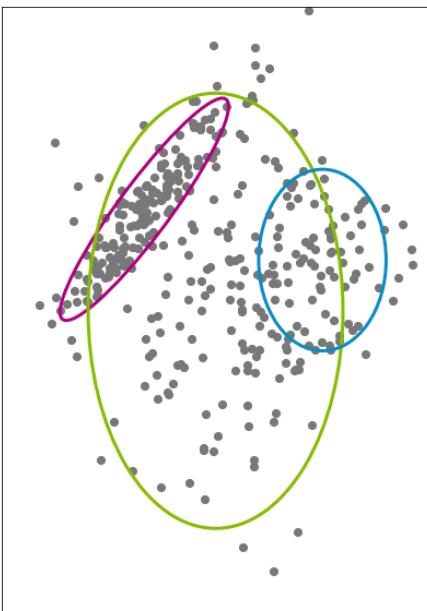
$$\text{p}(\mathbf{x} | \mu_k, \Sigma_k, \pi_k)$$

E step: estimate responsibilities



Compute $r_{ik} = \Pr(z_{ik} = 1 | x_i, \text{params}\{\boldsymbol{\pi}_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\})$
(responsibilities)

Responsibilities in equations



$$\Pr(z_{ik} = 1 \mid x_i, \theta)$$

Responsibility cluster k takes for observation i

$$r_{ik} = \frac{\pi_k N(x_i \mid \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_i \mid \mu_j, \Sigma_j)}$$

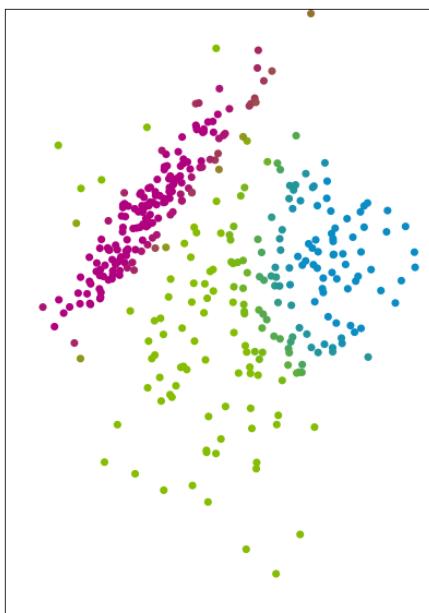
Normalized over all possible cluster assignments

M step:

Given the soft assignments r_{ij} , estimate parameters

Estimating cluster parameters from soft assignments

$$r_{ij}$$



Estimating cluster parameters from assignments r_{ij}

R	G	B	Cluster
$x_1[1]$	$x_1[2]$	$x_1[3]$	3
$x_2[1]$	$x_2[2]$	$x_2[3]$	3
$x_3[1]$	$x_3[2]$	$x_3[3]$	3
$x_4[1]$	$x_4[2]$	$x_4[3]$	1
$x_5[1]$	$x_5[2]$	$x_5[3]$	2
$x_6[1]$	$x_6[2]$	$x_6[3]$	2

Suppose that magically r_{ik} 's are all 0, 1, i.e. hard assignments

Estimate $\{\pi_k, \mu_k, \Sigma_k\}$ given data assigned to cluster k

Mean/covariance MLE

R	G	B	Cluster
x ₁ [1]	x ₁ [2]	x ₁ [3]	3
x ₂ [1]	x ₂ [2]	x ₂ [3]	3
x ₃ [1]	x ₃ [2]	x ₃ [3]	3

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{i \text{ in } k} x_i$$

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{i \text{ in } k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

Cluster proportion MLE

R	G	B	Cluster
x ₄ [1]	x ₄ [2]	x ₄ [3]	1

R	G	B	Cluster
x ₅ [1]	x ₅ [2]	x ₅ [3]	2
x ₆ [1]	x ₆ [2]	x ₆ [3]	2

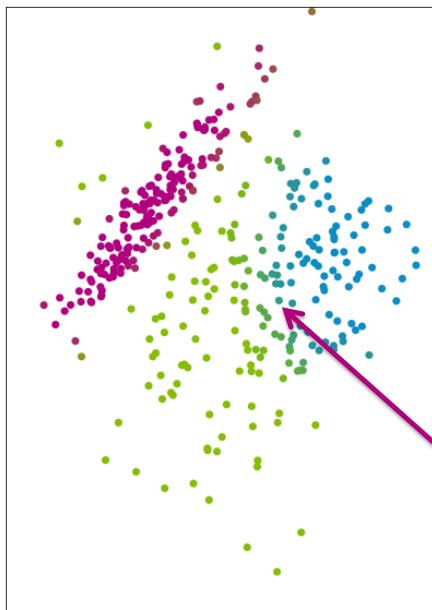
R	G	B	Cluster
x ₁ [1]	x ₁ [2]	x ₁ [3]	3
x ₂ [1]	x ₂ [2]	x ₂ [3]	3
x ₃ [1]	x ₃ [2]	x ₃ [3]	3

obs in cluster k

$$\hat{\pi}_k = \frac{N_k}{N}$$

total # of obs

Estimating cluster parameters from soft assignments



Instead of having a full observation x_i in cluster k , just allocate a portion r_{ik}

x_i divided across all clusters,
as determined by r_{ik}

Maximum likelihood estimation from soft assignments

R	G	B	r_{i1}	r_{i2}	r_{i3}
$x_1[1]$	$x_1[2]$	$x_1[3]$	0.30	0.18	0.52
$x_2[1]$	$x_2[2]$	$x_2[3]$	0.01	0.26	0.73
$x_3[1]$	$x_3[2]$	$x_3[3]$	0.002	0.008	0.99
$x_4[1]$	$x_4[2]$	$x_4[3]$	0.75	0.10	0.15
$x_5[1]$	$x_5[2]$	$x_5[3]$	0.05	0.93	0.02
$x_6[1]$	$x_6[2]$	$x_6[3]$	0.13	0.86	0.01

Total weight in cluster:
(effective # of obs)

1.242 2.8 2.42

52% chance this obs is in cluster 3

Maximum likelihood estimation from soft assignments

R	G	B	Cluster 1 weights	
x ₁ [1]	x ₁ [2]	x ₁ [3]	0.30	
x ₂ [1]	R	G	B	Cluster 2 weights
x ₃ [1]				
x ₄ [1]	x ₁ [1]	x ₁ [2]	x ₁ [3]	0.18
x ₅ [1]	x ₂ [1]	R	G	B
x ₆ [1]	x ₃ [1]			Cluster 3 weights
x ₄ [1]	x ₁ [1]	x ₁ [2]	x ₁ [3]	0.52
x ₅ [1]	x ₂ [1]	x ₂ [2]	x ₂ [3]	0.73
x ₆ [1]	x ₃ [1]	x ₃ [2]	x ₃ [3]	0.99
	x ₄ [1]	x ₄ [2]	x ₄ [3]	0.15
	x ₅ [1]	x ₅ [2]	x ₅ [3]	0.02
	x ₆ [1]	x ₆ [2]	x ₆ [3]	0.01



Cluster-specific location/shape MLE

R	G	B	Cluster 1 weights
x ₁ [1]	x ₁ [2]	x ₁ [3]	0.30
x ₂ [1]	x ₂ [2]	x ₂ [3]	0.01
x ₃ [1]	x ₃ [2]	x ₃ [3]	0.002
x ₄ [1]	x ₄ [2]	x ₄ [3]	0.75
x ₅ [1]	x ₅ [2]	x ₅ [3]	0.05
x ₆ [1]	x ₆ [2]	x ₆ [3]	0.13

1.242

Compute cluster parameter estimates
with weights on each row operation

$$\hat{\mu}_k = \frac{1}{N_k^{\text{soft}}} \sum_{i=1}^N r_{ik} x_i$$

$$\hat{\Sigma}_k = \frac{1}{N_k^{\text{soft}}} \sum_{i=1}^N r_{ik} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

$$N_k^{\text{soft}} = \sum_{i=1}^N r_{ik}$$

Total weight in cluster k
= effective # obs

MLE of cluster proportions $\hat{\pi}_k$

r_{i1}	r_{i2}	r_{i3}
0.30	0.18	0.52
0.01	0.26	0.73
0.002	0.008	0.99
0.75	0.10	0.15
0.05	0.93	0.02
0.13	0.86	0.01

Total weight in cluster:



Total weight in dataset:

6

datapoints N

©2017 Emily Fox

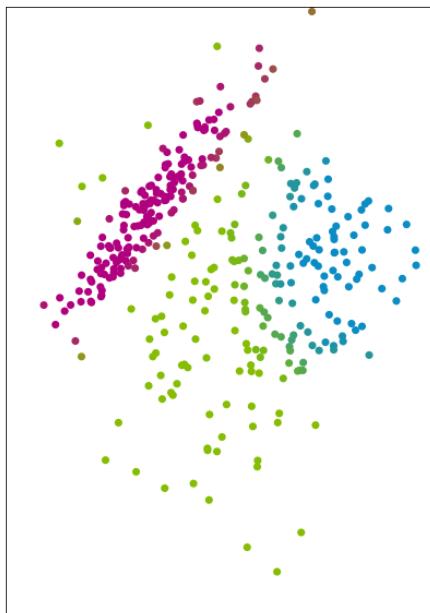
$$\hat{\pi}_k = \frac{N_k^{\text{soft}}}{N}$$

$$N_k^{\text{soft}} = \sum_{i=1}^N r_{ik}$$

Total weight in cluster k
= effective # obs

Estimate cluster proportions from relative weights

M step summary



Compute cluster parameter estimates from soft assignments

$$\hat{\mu}_k = \frac{1}{N_k^{\text{soft}}} \sum_{i=1}^N r_{ik} x_i$$

$$\hat{\Sigma}_k = \frac{1}{N_k^{\text{soft}}} \sum_{i=1}^N r_{ik} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

$$N_k^{\text{soft}} = \sum_{i=1}^N r_{ik} \quad \hat{\pi}_k = \frac{N_k^{\text{soft}}}{N}$$

Expectation maximization (EM)

Expectation maximization (EM): An iterative algorithm

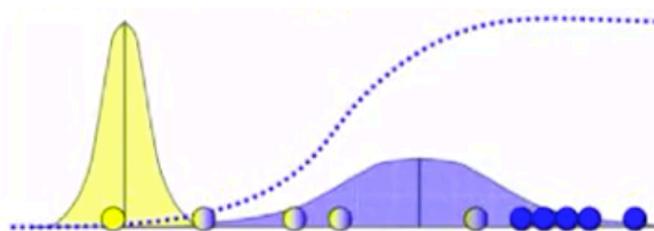
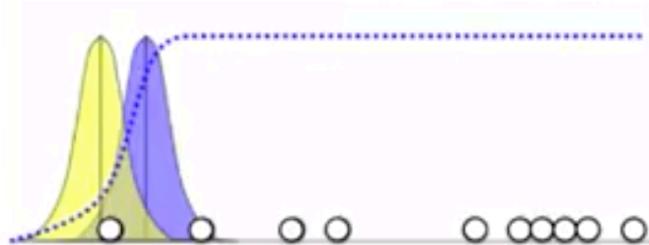
Motivates an iterative algorithm: $\Theta = \{\pi^{(t)}, \mu^{(t)}, \Sigma^{(t)}\}$

1. E-step: estimate cluster responsibilities given current parameter estimates

$$\hat{r}_{ik} = \frac{\hat{\pi}_k N(x_i | \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{j=1}^K \hat{\pi}_j N(x_i | \hat{\mu}_j, \hat{\Sigma}_j)}$$

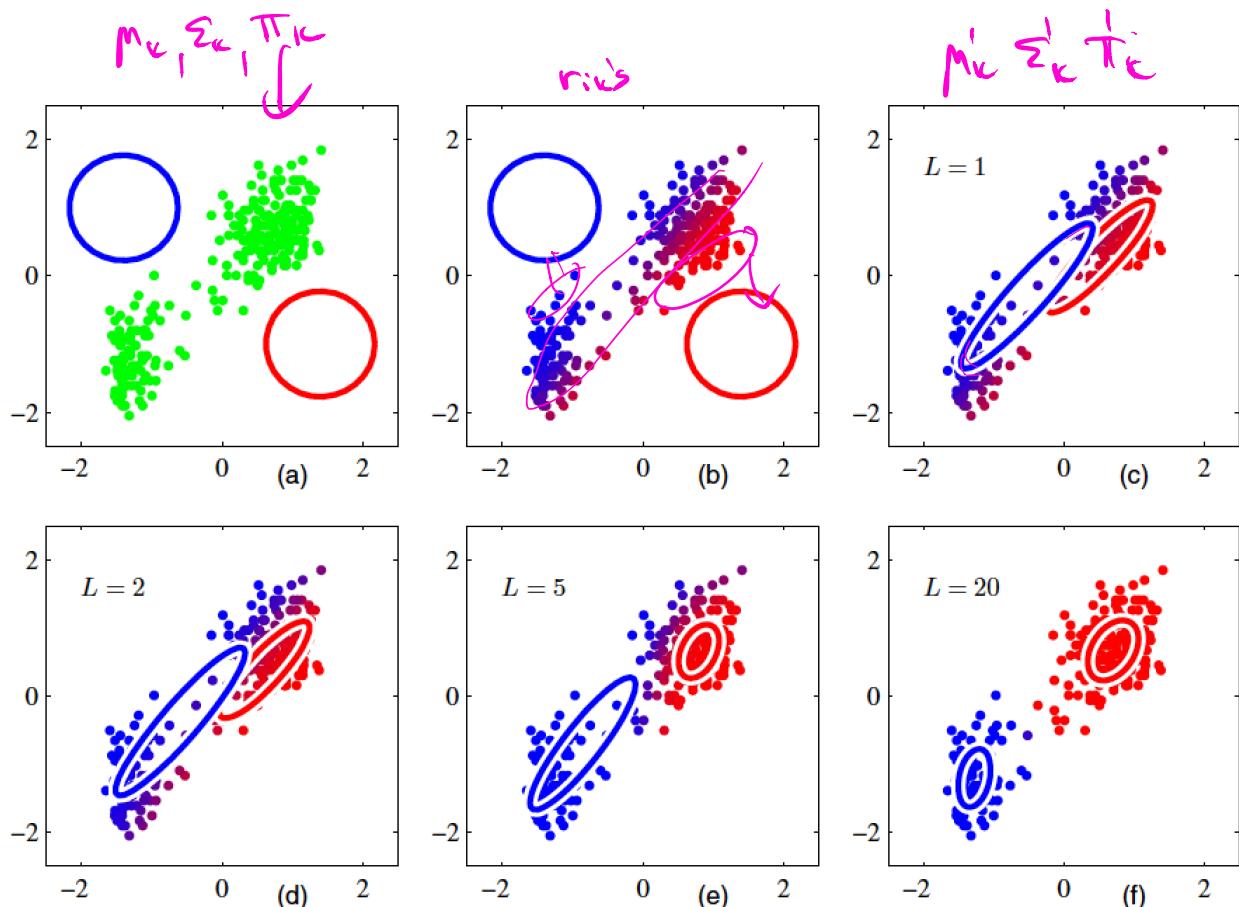
2. M-step: maximize likelihood over parameters given current responsibilities

$$\hat{\pi}_k, \hat{\mu}_k, \hat{\Sigma}_k | \{\hat{r}_{ik}, x_i\}$$

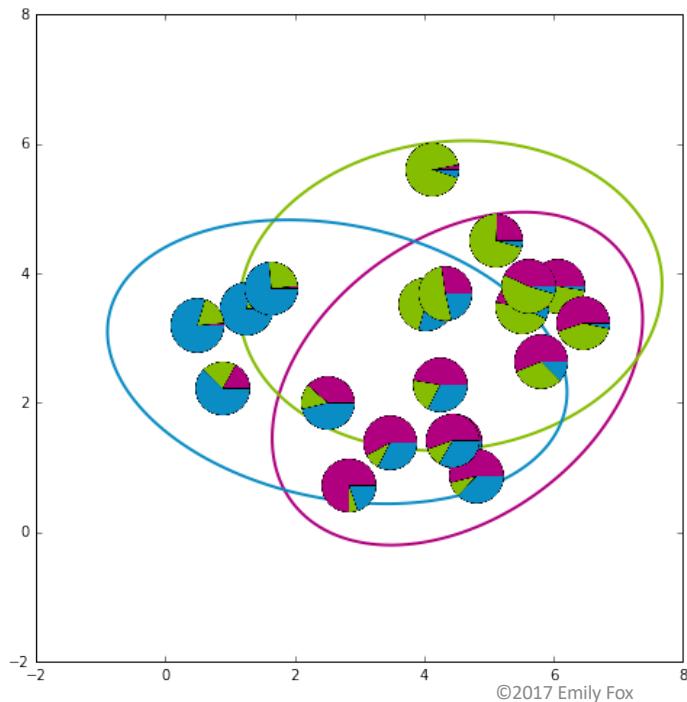


©2017 Emily Fox

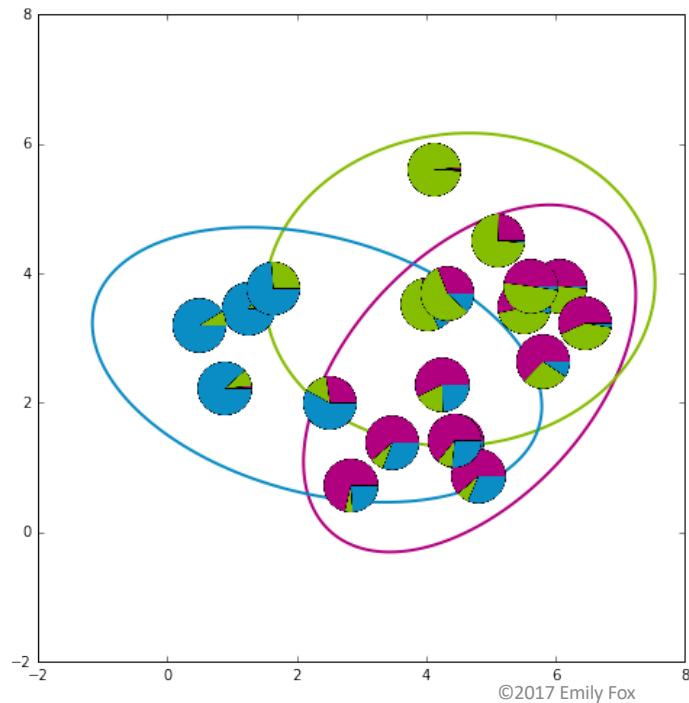
CSE 446: Machine Learning



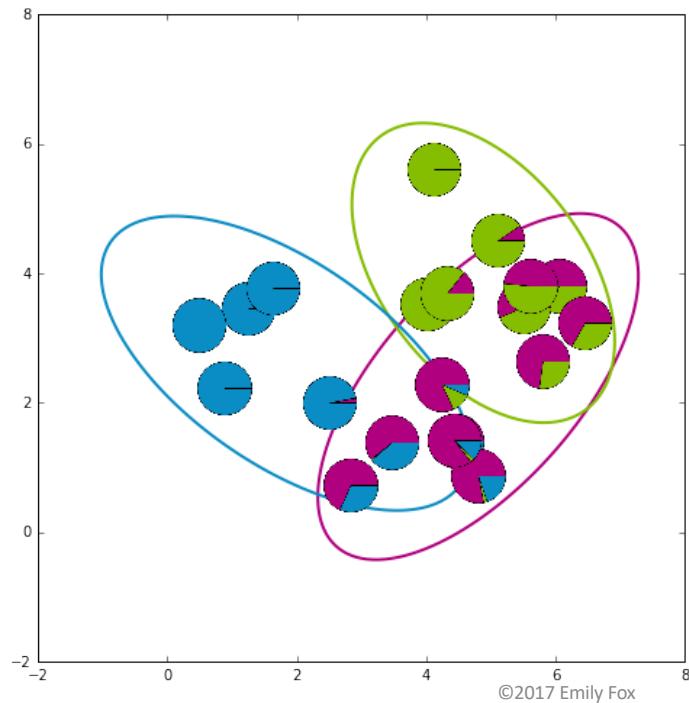
EM for mixtures of Gaussians in pictures – initialization



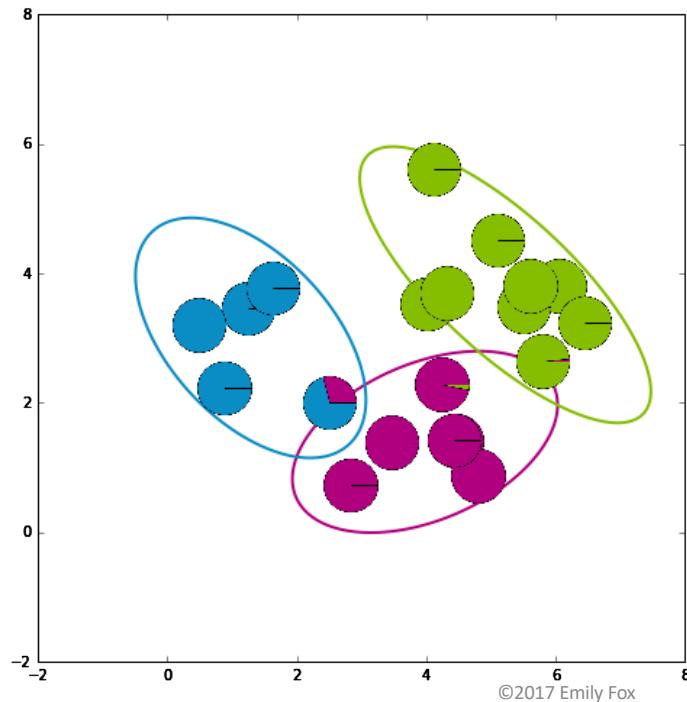
EM for mixtures of Gaussians in pictures – after 1st iteration



EM for mixtures of Gaussians in pictures – after 2nd iteration



EM for mixtures of Gaussians in pictures – converged solution



The nitty gritty of EM

Convergence and initialization of EM

Convergence of EM

- EM is a **coordinate-ascent algorithm**
 - Can equate E-and M-steps with alternating maximizations of the objective function
- Convergence to a **local maximum.**
- We assess via (log) likelihood of data under current parameter and responsibility estimates

Initialization

- Many ways to initialize the EM algorithm
- Important for convergence rates & quality of local maximum found
- Examples:
 - Choose K observations at random to define K “centroids”. Assign other observations to nearest centroid to form initial parameter estimates.
 - Initialize from k-means solution

Potential of vanilla EM to overfit

Overfitting of MLE

Maximizing likelihood can **overfit to data**

Imagine at K=2 example with one obs assigned to **cluster 1** and others assigned to **cluster 2**

- What parameter values maximize likelihood?



Set center equal to point and shrink variance to 0

Likelihood goes to ∞ !

Simple regularization of M-step for mixtures of Gaussians

Simple fix: **Don't let variances $\rightarrow 0$!**

Add small amount to diagonal of covariance estimate

Summary

What you can do now...

- Understand Gaussian mixture models (and multivariate Gaussians)
- Estimate soft assignments (responsibilities) given mixture model parameters
- Solve maximum likelihood parameter estimation using soft assignments (weighted data)
- Implement an EM algorithm for inferring soft assignments and cluster parameters
 - Determine an initialization strategy
 - Implement a variant that helps avoid overfitting issues
- Compare and contrast with k-means
 - Soft vs. hard assignments