

Mixture Models and EM: Model-Based Clustering

CSE 446: Machine Learning

Slides by Emily Fox

University of Washington

May 22, 2019

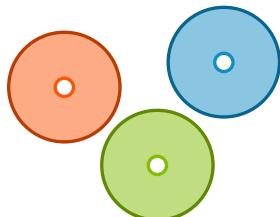
Limitations of k-means

Assigns observations to closest cluster center

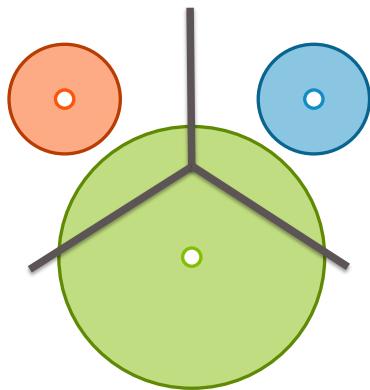
$$z_i \leftarrow \arg \min_j \|\mu_j - \mathbf{x}_i\|_2^2$$

Only center matters
Not cluster shapes

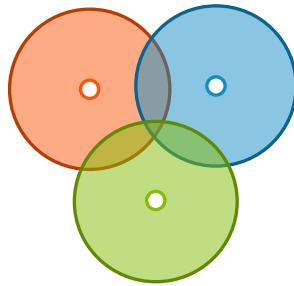
Equivalent to assuming
spherically symmetric clusters



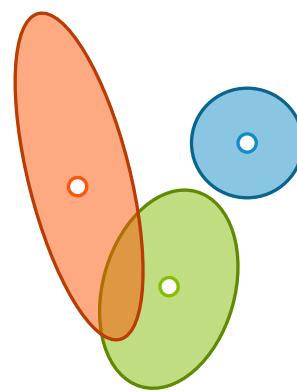
Failure modes of k-means



disparate cluster sizes



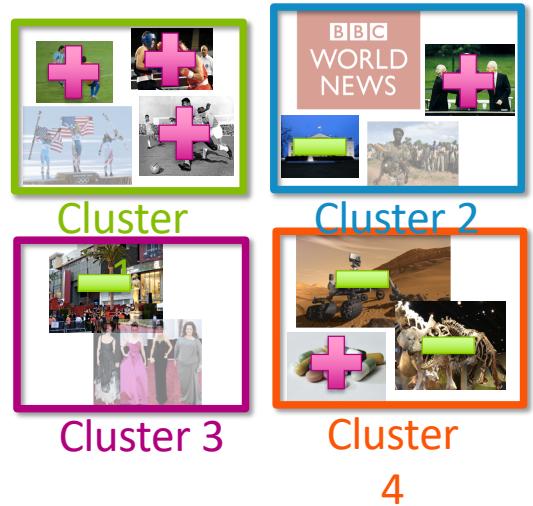
overlapping clusters



different
shaped/oriented
clusters

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**

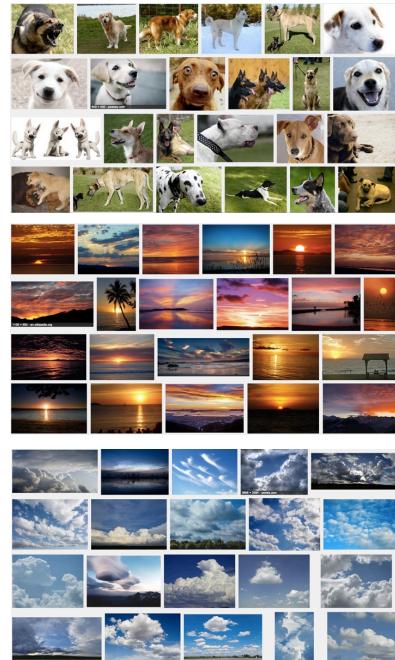
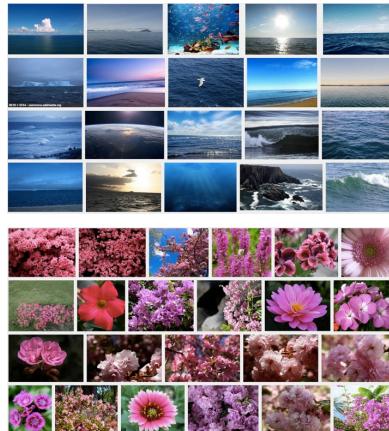


Mixture models

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

average RGB value over all pixels in entire image.



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



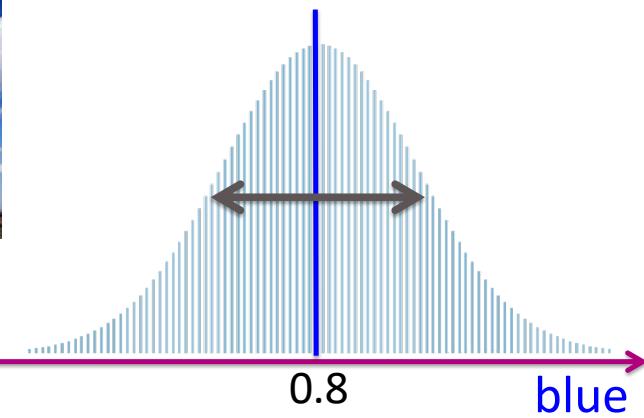
[R = 0.85, G = 0.05, B = 0.35]



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

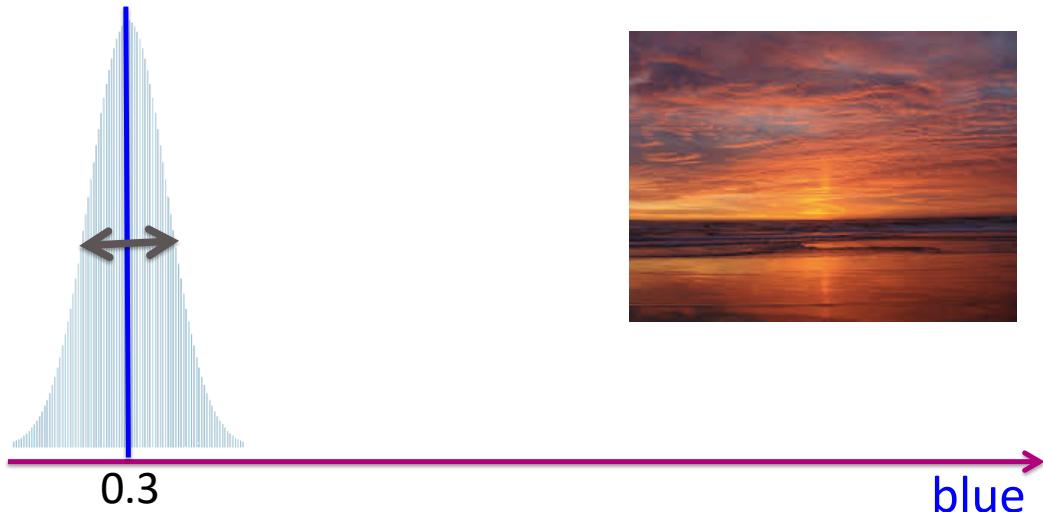
Distribution over all **cloud** images

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



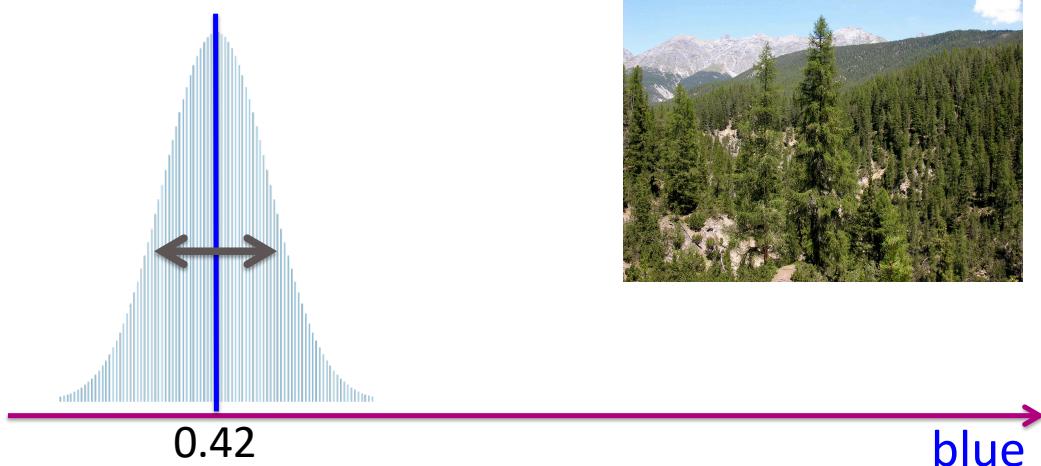
Distribution over all sunset images

Let's look at just the blue dimension

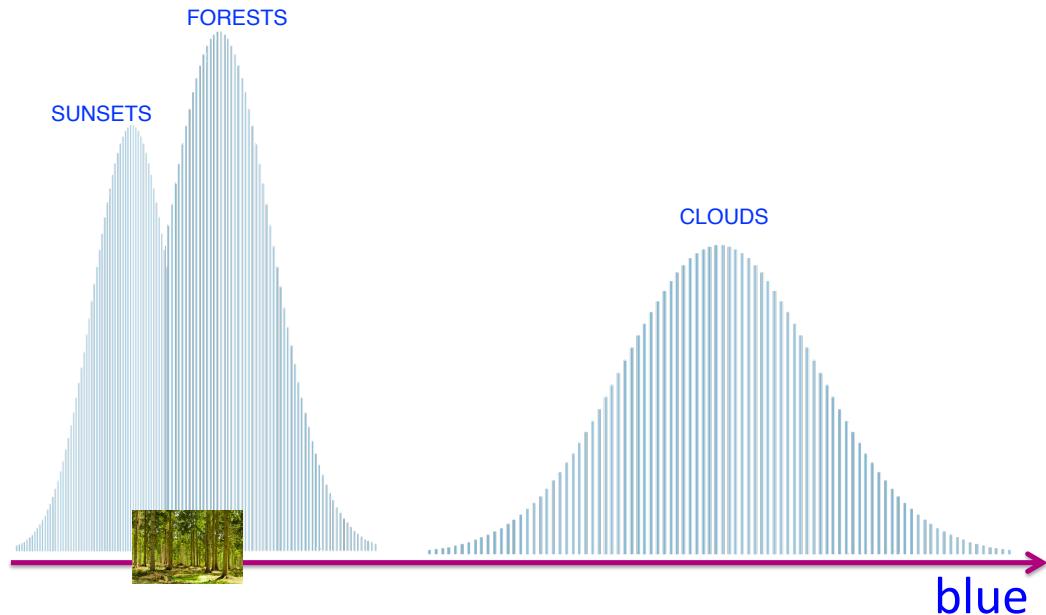


Distribution over all forest images

Let's look at just the blue dimension



Distribution over **all** images



Can be distinguished along other dim

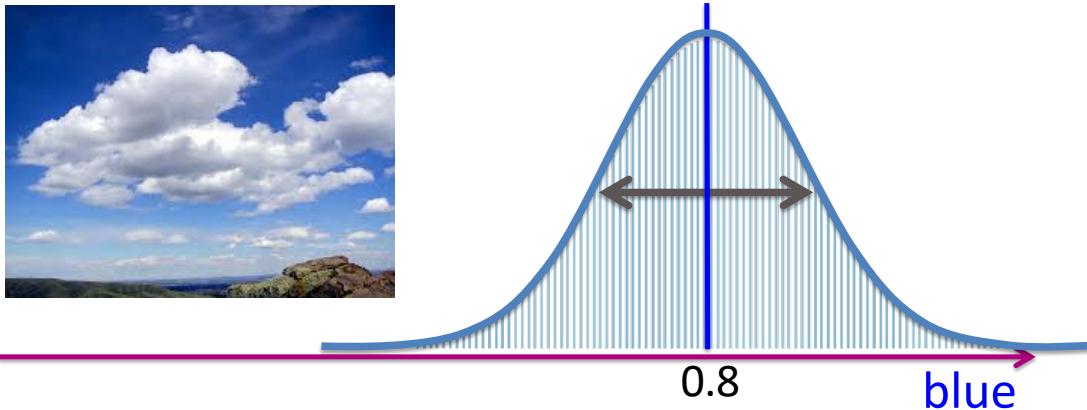
Now look at the **red** dimension



Background: Gaussian distributions

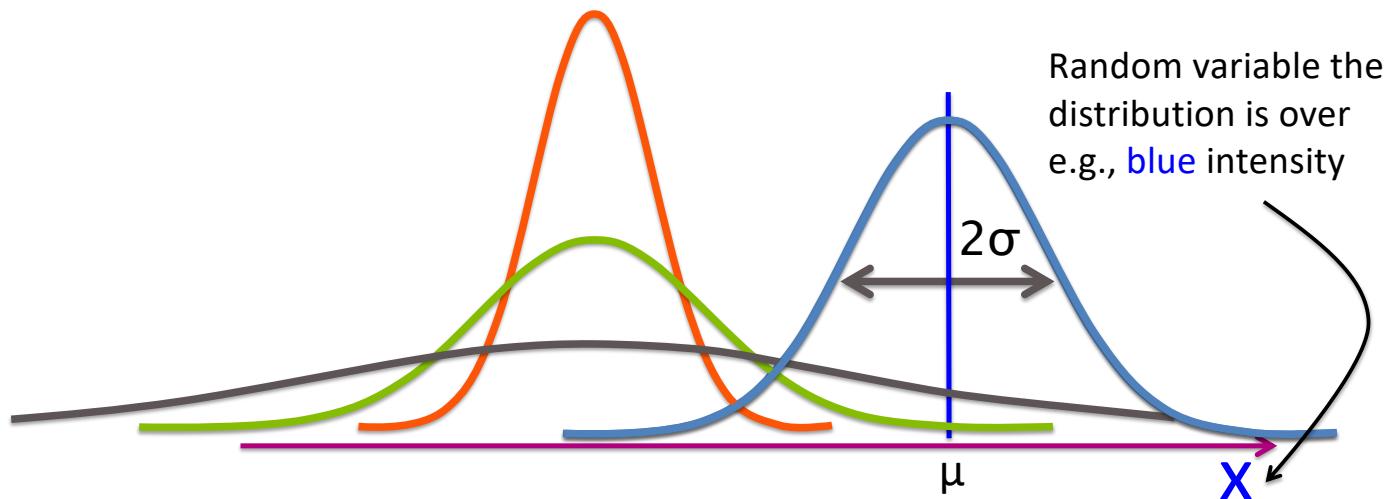
Model for a given image type

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



1D Gaussians

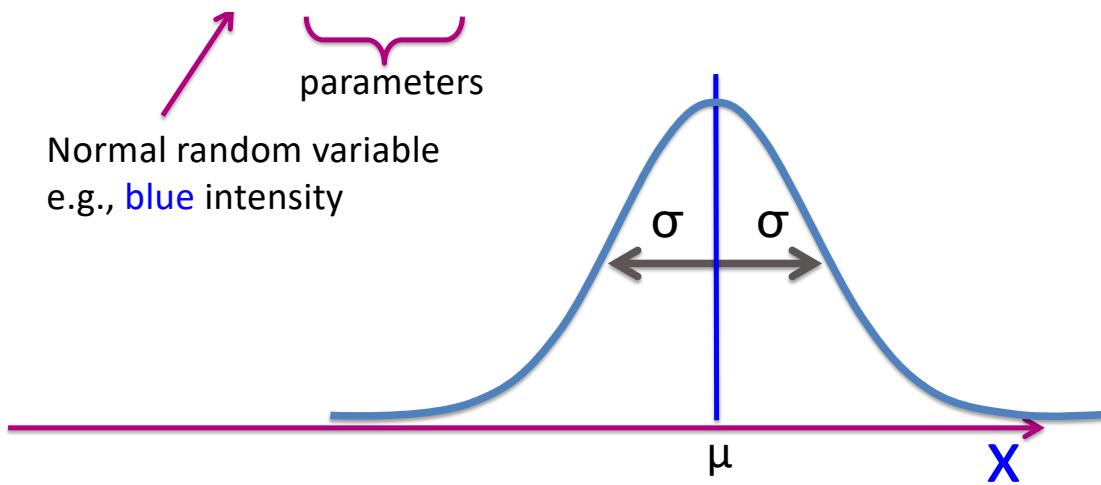
Fully specified by mean μ and variance σ^2 (or st. dev. σ)



Notation a 1D Gaussian distribution

Probability that X takes on the value x .

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



Multivariate Gaussian density

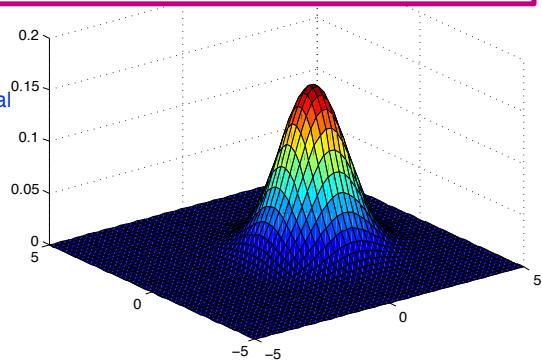
$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma)$

parameters

Random vector
e.g., [R, G, B] intensities

determinant
 $= (2\pi)^d |\Sigma|$



quadratic form

Note that covariance matrix is always positive definite
(unless degenerate; i.e. one of the variables has no variance)

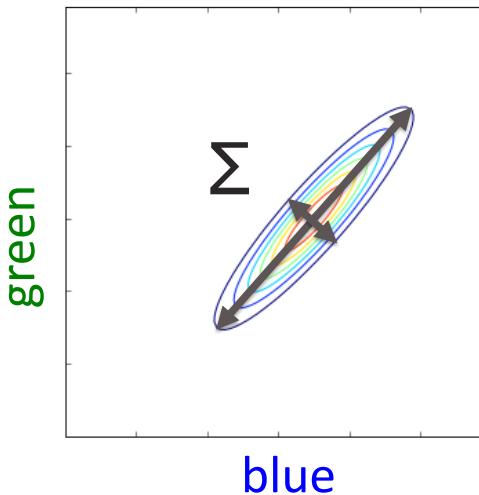
Multivariate Gaussian

Fully specified by mean μ and covariance Σ

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

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

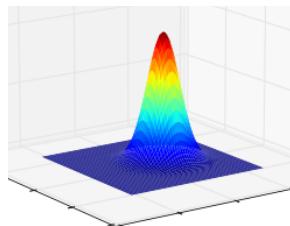
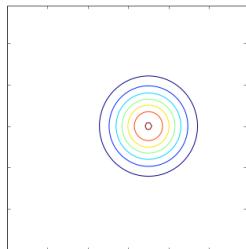
covariance determines
orientation + spread



Spread in this direction is σ_{blue}^2

Covariance structure

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

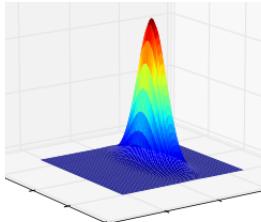
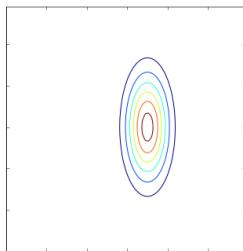


$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

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

Covariance structure

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



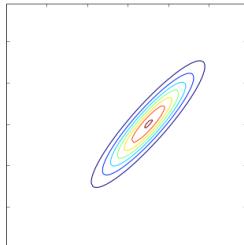
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu} = 0, \Sigma) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} \right\}$$



$x_1/\sigma_1^2 + x_2/\sigma_2^2 = C \rightarrow \text{level sets are ellipses!}$

Covariance Structure

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



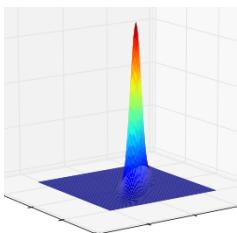
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu} = 0, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right\}$$

eigenvector decomposition
↓
 $\Sigma = QDQ^T \rightarrow \Sigma^{-1} = QD^{-1}Q^T$

$$\mathbf{x}^T Q D^{-1} Q^T \mathbf{x}$$

Let $\mathbf{y} = Q^T \mathbf{x}$ → projection of \mathbf{x} into the basis provided by Q

$$\begin{aligned} &= \mathbf{y}^T D^{-1} \mathbf{y} \\ &= y_1 / \sigma_1^2 + y_2 / \sigma_2^2 \rightarrow \text{ellipse in the new coordinate system specified by columns of } Q \\ &\quad (\text{the eigenvectors}) \end{aligned}$$



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

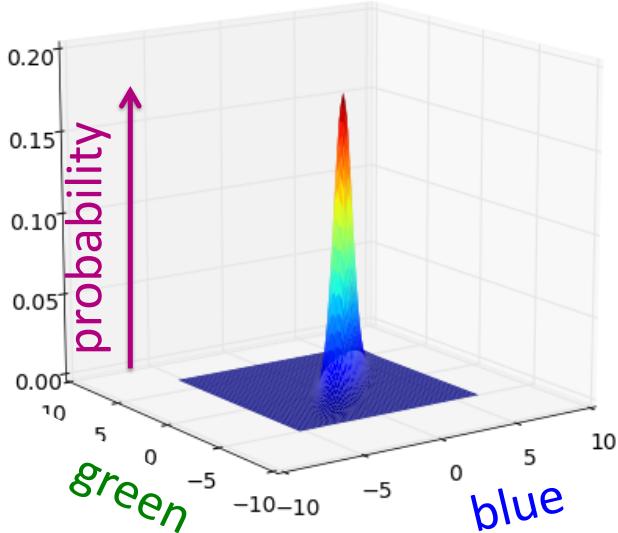
1. Start with $\mathbf{X} \sim N(0, I)$
2. (Scaling step) $D^{1/2}\mathbf{X} \sim N(0, D)$.
3. (Rotation) $UD^{1/2}\mathbf{X} \sim N(0, \Sigma)$ where $\Sigma = UDU^T$.
4. (Translation) $UD^{1/2}\mathbf{X} + \boldsymbol{\mu} \sim N(\boldsymbol{\mu}, \Sigma)$



Multivariate Gaussians

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

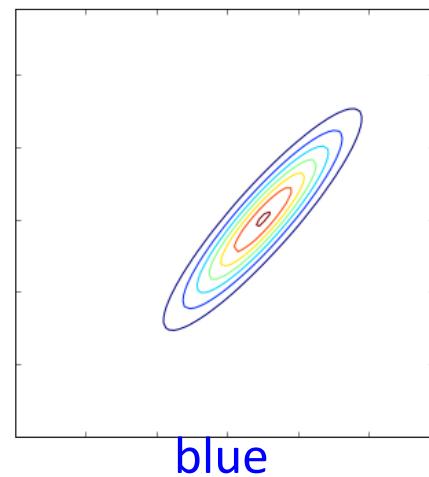
3D mesh plot



24

©2017 Emily Fox

Contour plot



CSE 446: Machine Learning

Summary

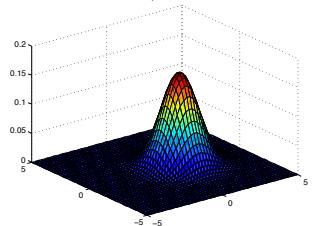
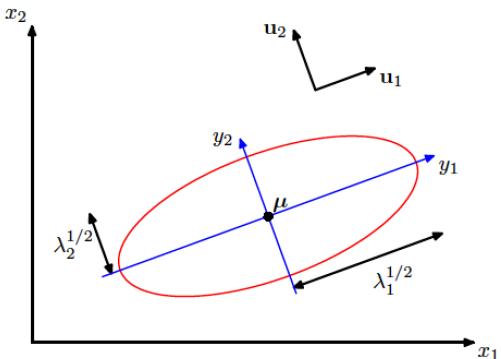
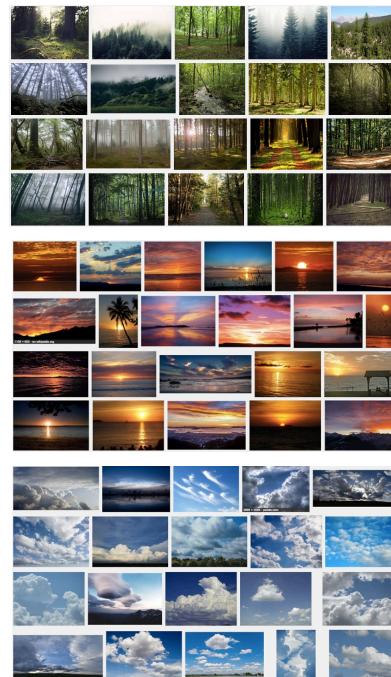
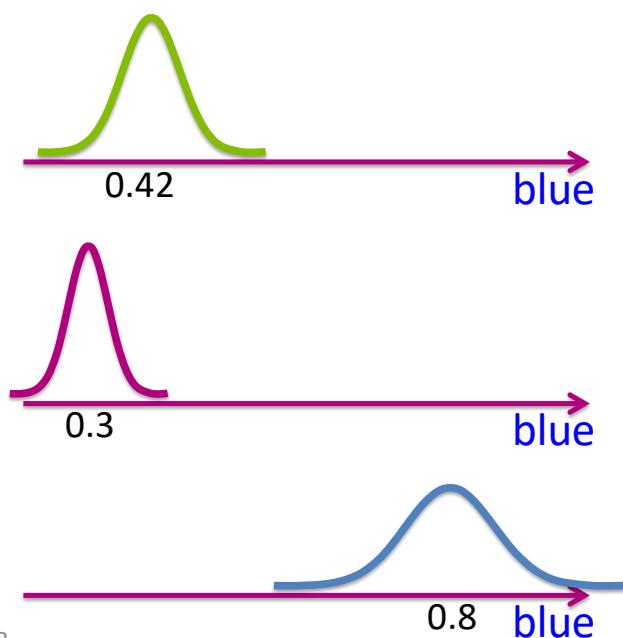


Figure 2.7 The red curve shows the elliptical surface of constant probability density for a Gaussian in a two-dimensional space $x = (x_1, x_2)$ on which the density is $\exp(-1/2)$ of its value at $x = \mu$. The major axes of the ellipse are defined by the eigenvectors u_i of the covariance matrix,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



Mixture Model (to be used for clustering)

Model as Gaussian per category/cluster



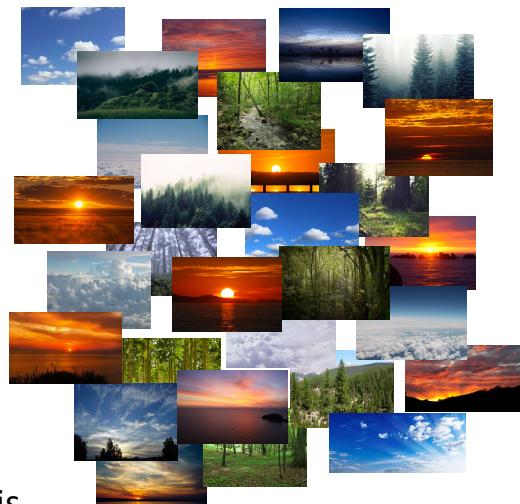
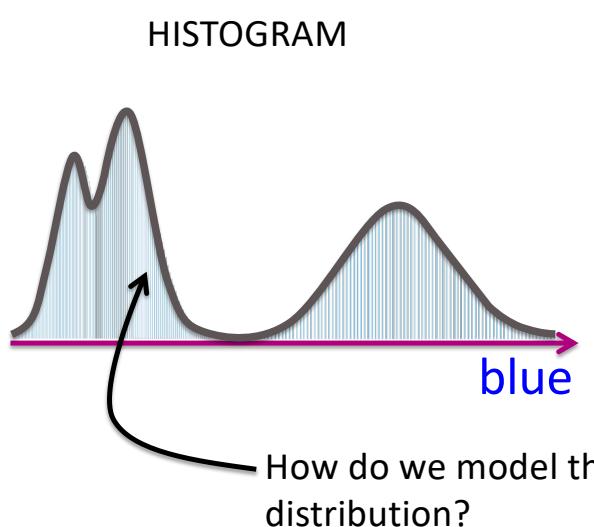
Forests

Sunsets

Clouds

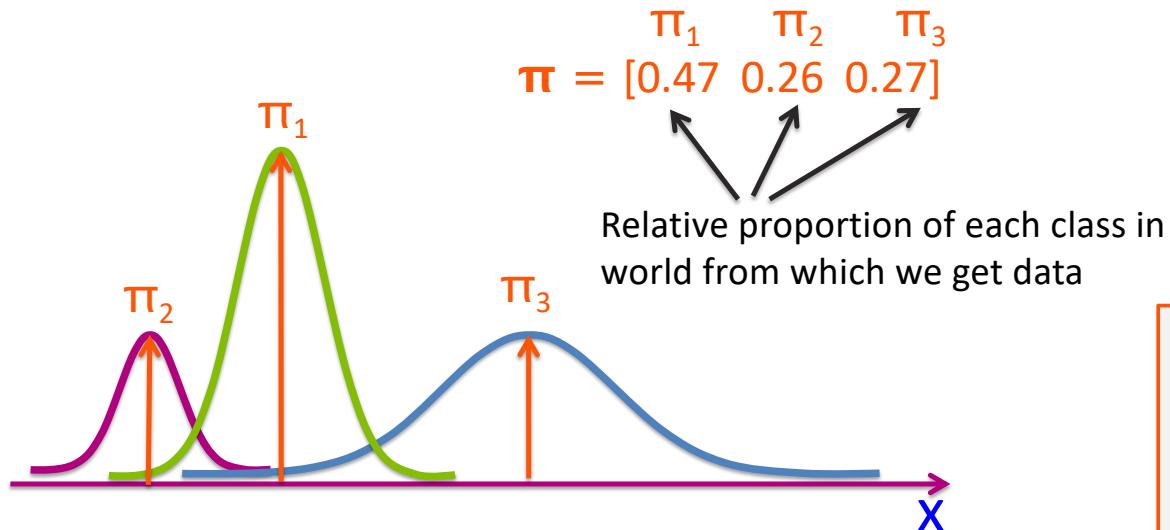
©2017 Emily Fox

Jumble of unlabeled images



Combination of weighted Gaussians

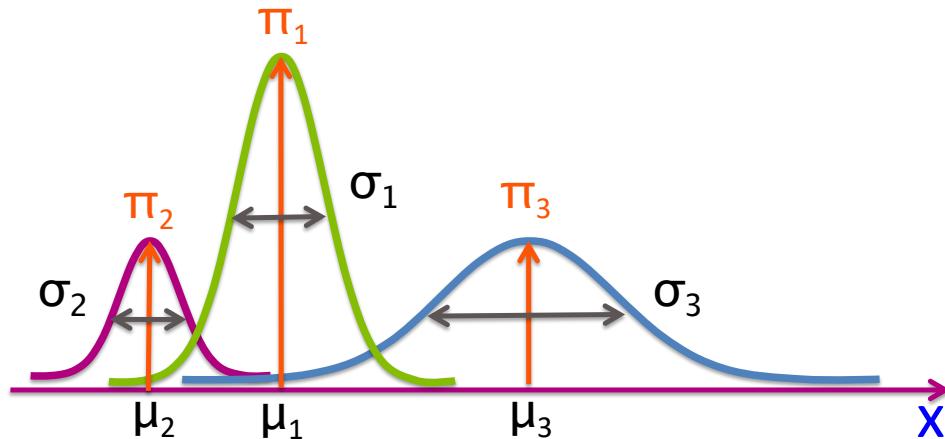
Associate a weight π_k with each Gaussian component



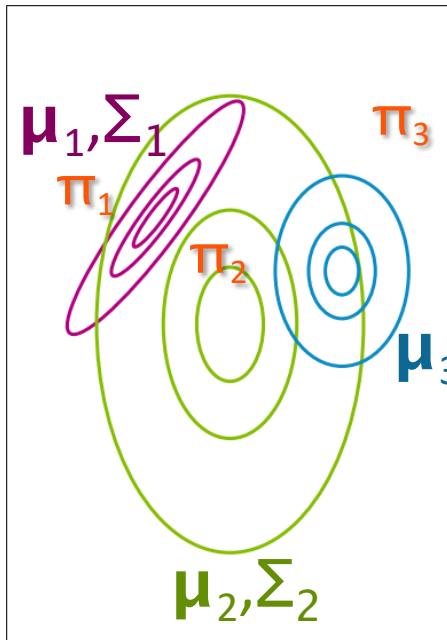
$$0 \leq \pi_k \leq 1$$
$$\sum_{k=1}^K \pi_k = 1$$

Mixture of Gaussians (1D)

Each mixture component represents a unique cluster specified by: $\{\pi_k, \mu_k, \sigma_k^2\}$



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

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

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

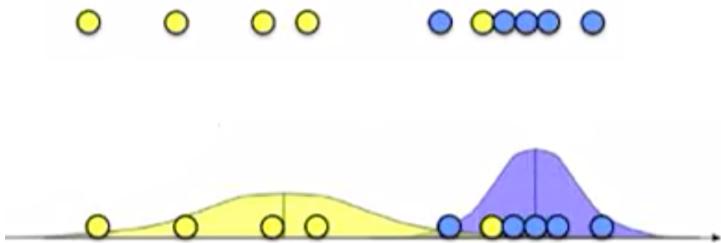


- Problem: Assume that the data comes from such a distribution, and recover the parameters of the distribution.
- Determine, for each point, the likelihood of it belonging to cluster j, for each j.

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

- Easy if know the source of each data point.

Could approximate the mean and variance for each of the K classes and compute the parameters from there.



- What if we don't know the source?



To understand better, introduce additional “latent” variables

- K clusters, defined by the following parameters

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



one hot encoding of the classes

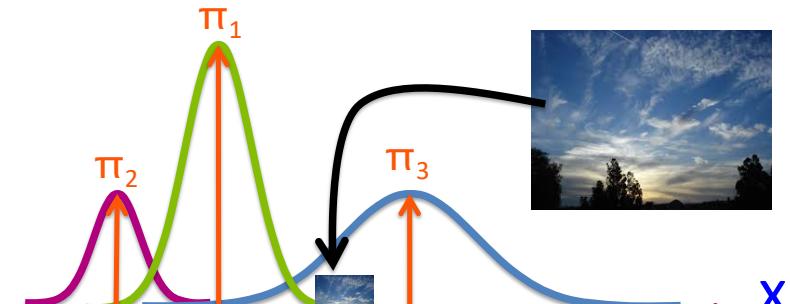
- For each point \mathbf{x} , let $\mathbf{z} = (z_1, \dots, z_K)$ indicate which cluster it was chosen from. These are called “latent variables”.

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

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

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

- Conditional distribution of x given z
- Therefore:
- Conditional probability of z given x



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 $p(z_k = 1|\mathbf{x}) = \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)}$

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



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.

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

MLE

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

Gradient = 0 gives the following conditions:

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \mathbf{x}_i$$

Probability that datapoint \mathbf{x}_i comes from the k th cluster

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

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T$$

$$\pi_k = \frac{N_k}{N}$$

N = total number of points in dataset

$$\begin{aligned} p(z_k = 1 | \mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \mu_j, \Sigma_j)}. \end{aligned}$$

Does not give us a closed form 😞

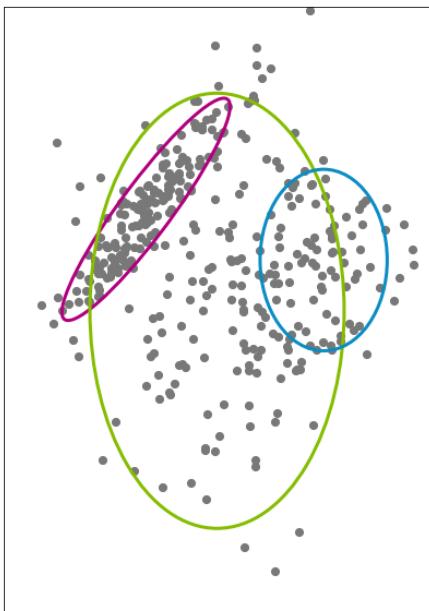
Expectation Maximization Algorithm

Algorithm for solving the Gaussian Mixture Model

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.

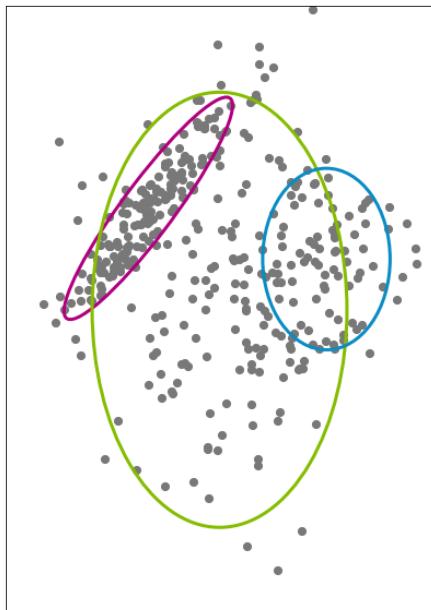
E step: estimate responsibilities



Using our current best guess of the parameters

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

Responsibilities in equations



Using whatever estimation of the parameters we currently have

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

Responsibility cluster k takes for observation i
Normalized over all possible cluster assignments

M step:

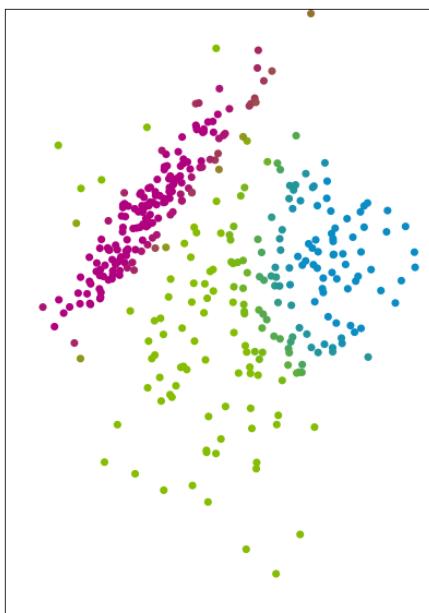
Given the soft assignments r_{ij} , estimate parameters

probability each point is in each cluster



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

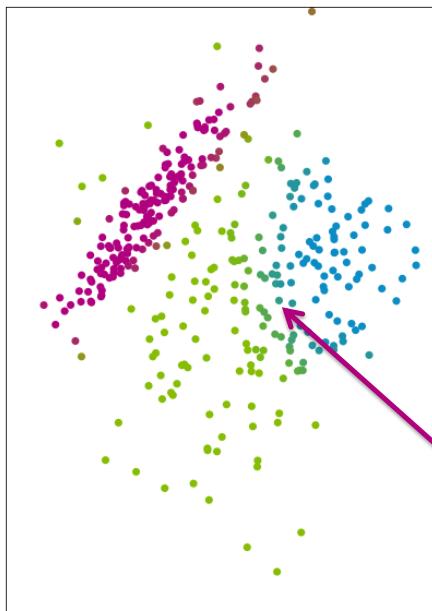
R	G	B	Cluster
x ₁ [1]	x ₁ [2]	x ₁ [3]	3
x ₂ [1]	x ₂ [2]	x ₂ [3]	3
x ₃ [1]	x ₃ [2]	x ₃ [3]	3

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

©2017 Emily Fox

expected number of points in BLUE cluster

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

Use these probabilities to update parameters of each class
 - estimate mean (weighted by prob point is in the cluster)
 - estimate variance (weighted by prob point is in cluster)

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:

1.242	2.8	2.42
-------	-----	------

Total weight in dataset:

6

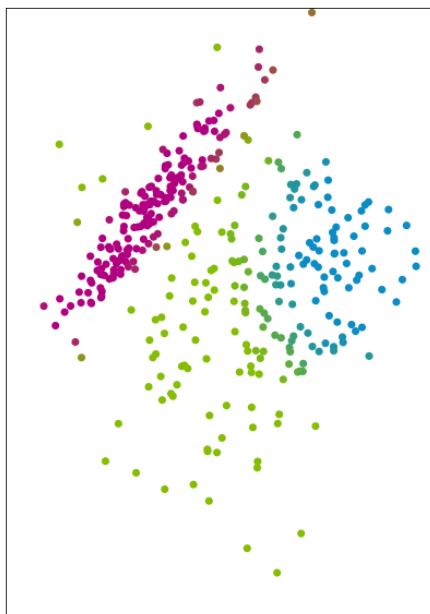
$$\hat{\pi}_k = \frac{N_k^{\text{soft}}}{N}$$

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

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

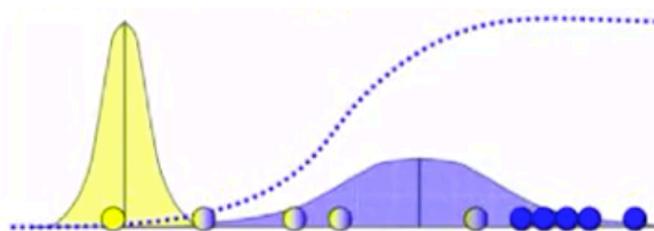
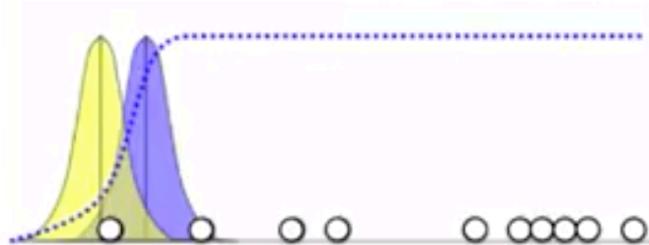
Motivates an iterative algorithm:

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

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

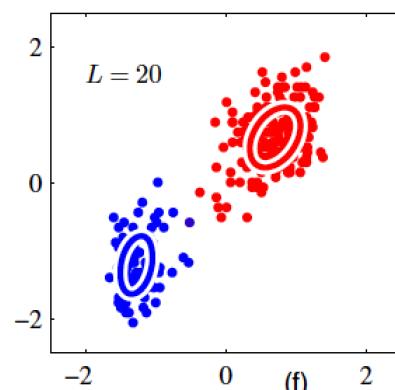
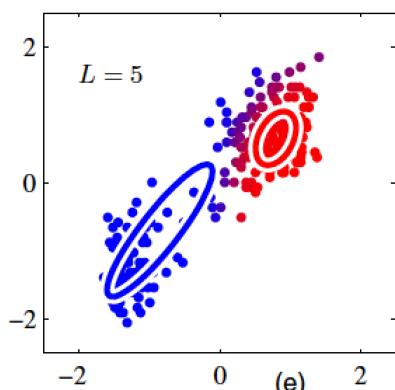
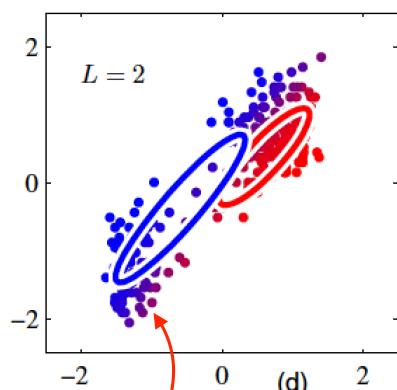
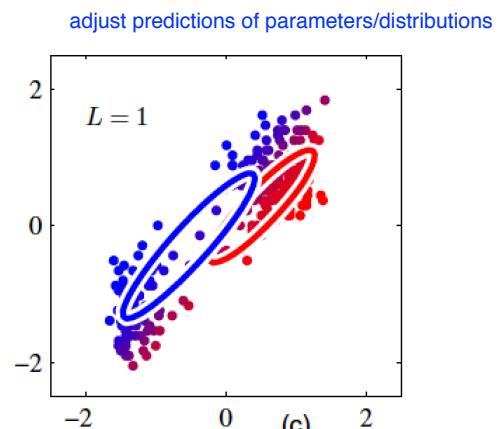
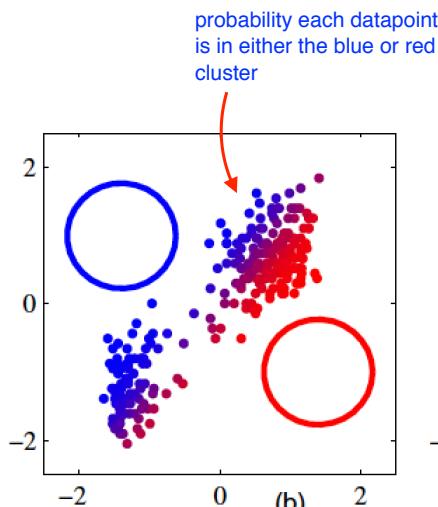
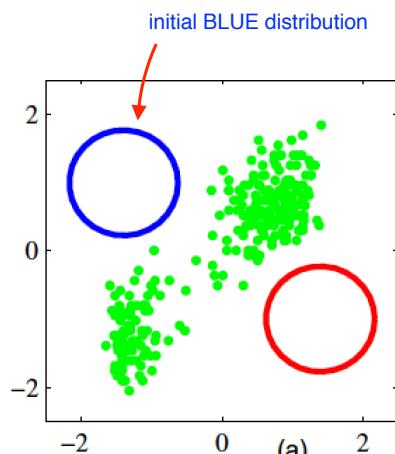
2. M-step: maximize likelihood over parameters given current responsibilities

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



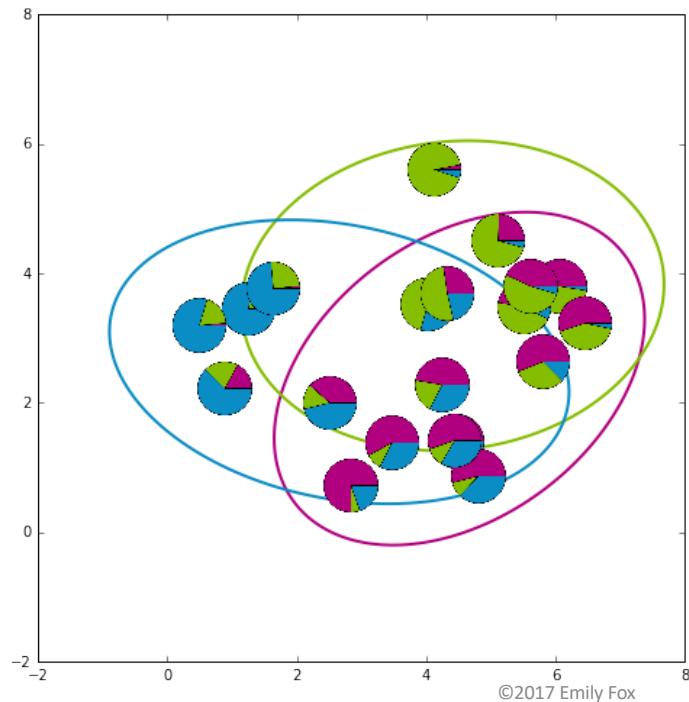
©2017 Emily Fox

CSE 446: Machine Learning

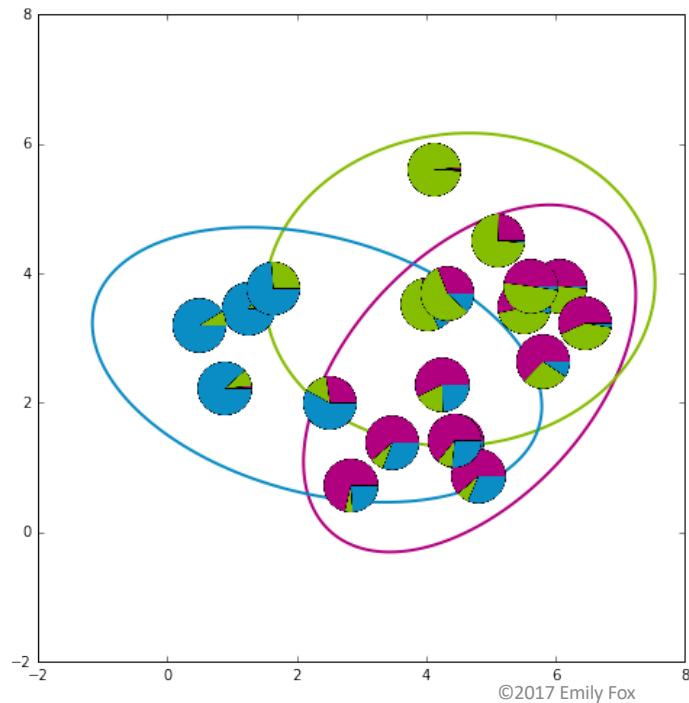


probability each datapoint
is in either the blue or red
cluster

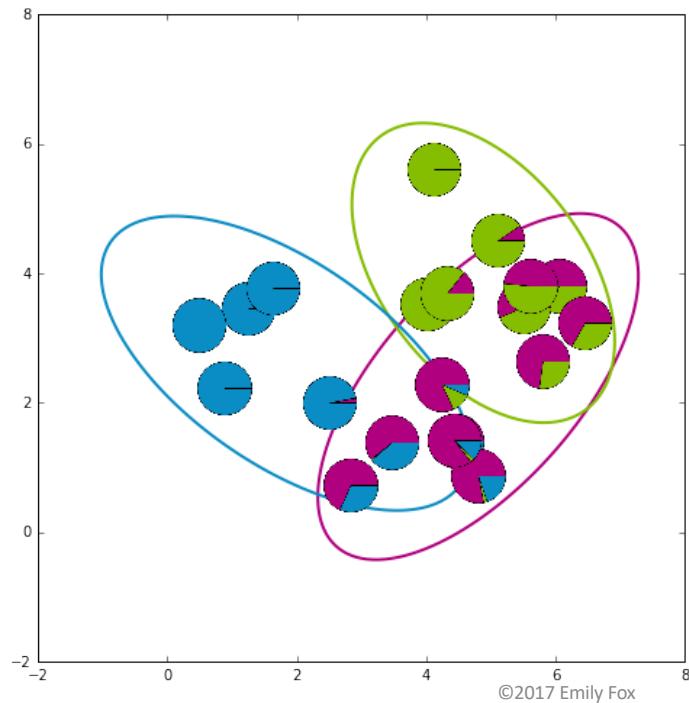
EM for mixtures of Gaussians in pictures – initialization



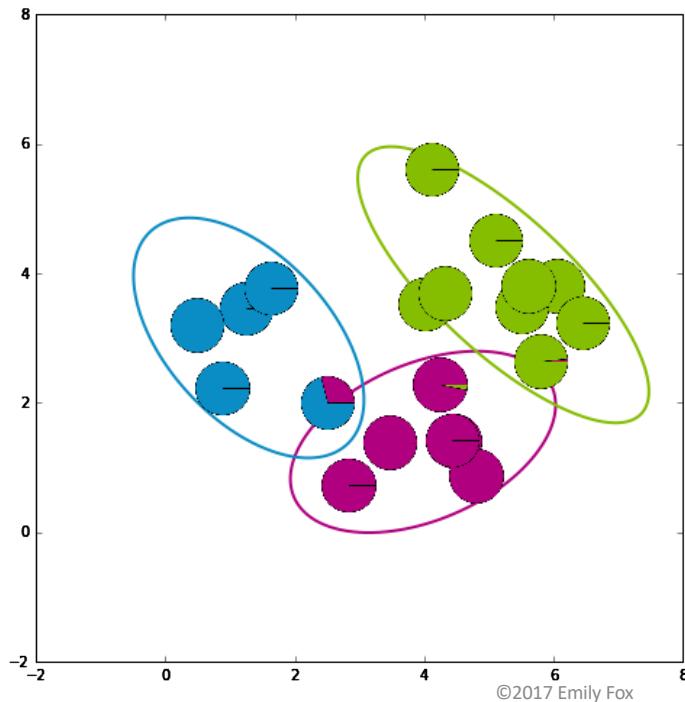
EM for mixtures of Gaussians in pictures – after 1st iteration



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

Not guaranteed to end up at a optimal solution; could get stuck in a local maximum

- 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

For convergence, check whether the distributions are moving

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

success of algorithm dependent on initial conditions; possibly use K means to determine initial locations of distributions

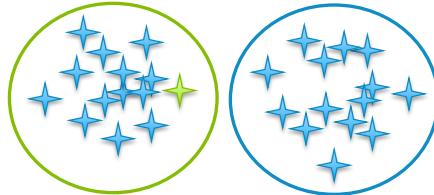
Potential of vanilla EM to overfit

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