

CSE 312

Foundations of Computing II

Lecture 27: Multivariate Gaussians, clustering and EM

Anna R. Karlin

Slide Credit: Emily Fox

Incorporating some of my own from CSE 446 ☺

Quiz 3 : Monday Dec 6
Finals : Monday Dec 13

Details on both at top of web page under announcements

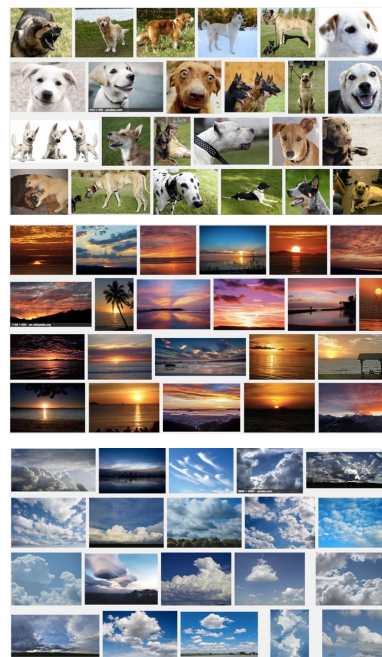
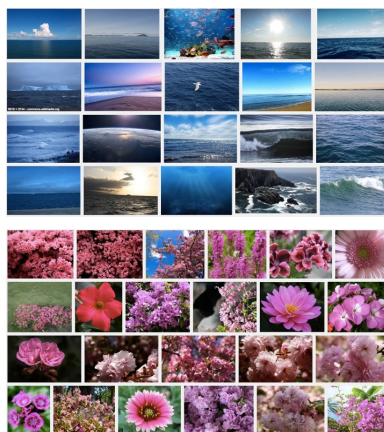
Goal for today

- Introduce you to a fundamental machine learning problem: clustering.
- Give you a very gentle introduction to multivariate Gaussian distributions.

Motivating application: Clustering images

Discover groups of similar images

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



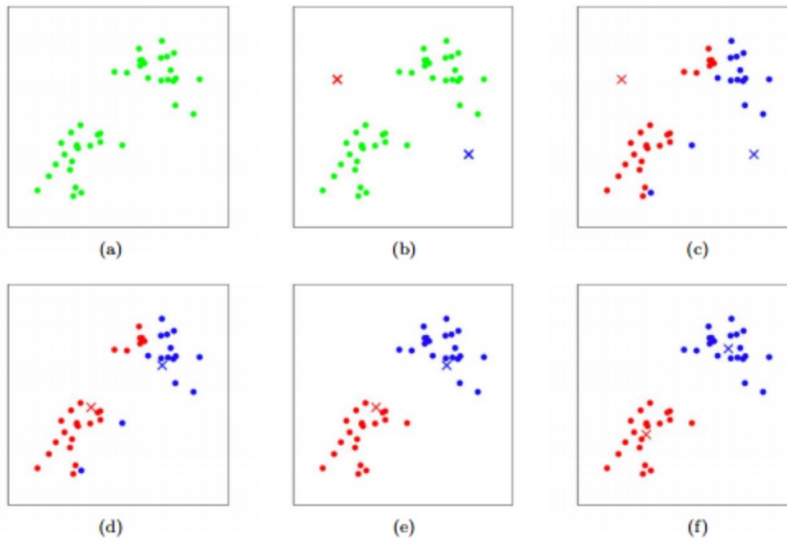
Another example: clustering documents

E.g. into international news, sports, culture, etc.

So how are these data represented?

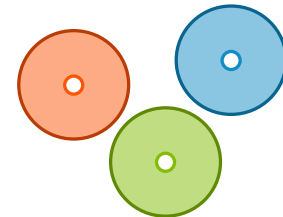
- As points in d -dimensional space (where d is typically large).

How to approach clustering? One way: k-means

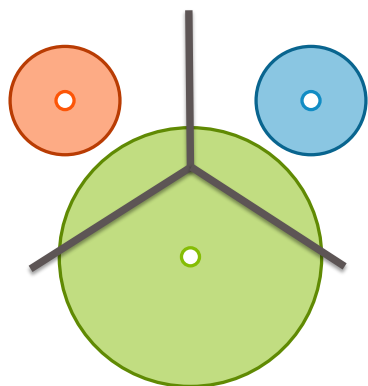


Only center matters
Not cluster shapes

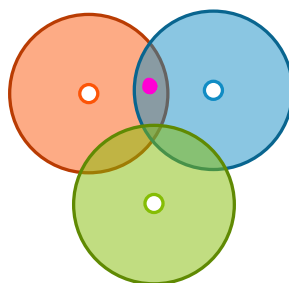
Equivalent to assuming
spherically symmetric clusters



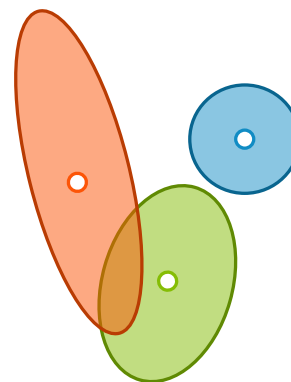
Failure modes of k-means clustering



disparate cluster sizes



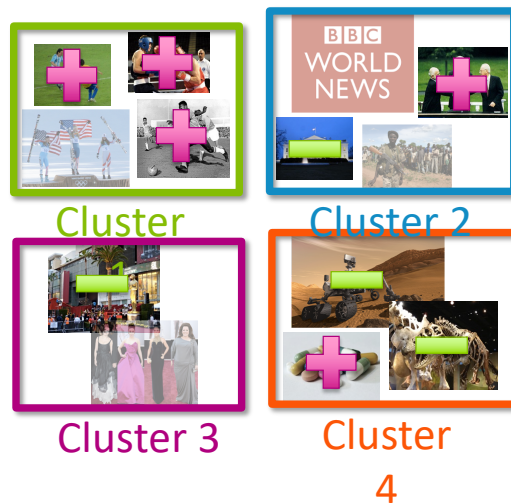
overlapping clusters



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
- Allow for cluster shapes not just centers
- Enables learning different weightings of dimensions
 - e.g., how much to weight each word in the vocabulary when computing cluster assignment



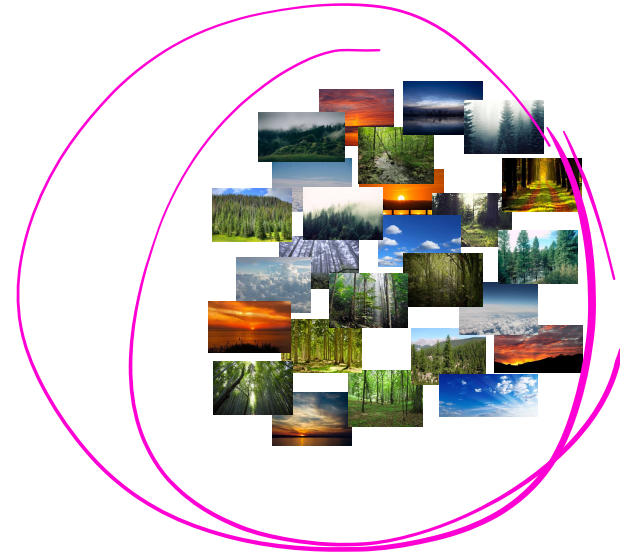


Mixture model

- k clusters, defined by probability distribution over **Gaussian** random variables.

- π_i, μ_i, σ_i^2 for each cluster. $\sum_{j=1}^k \pi_j = 1.$

- 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.





Background:
Multivariate Gaussian distributions

Overly simple image representation

Consider average red, green, blue pixel intensities



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



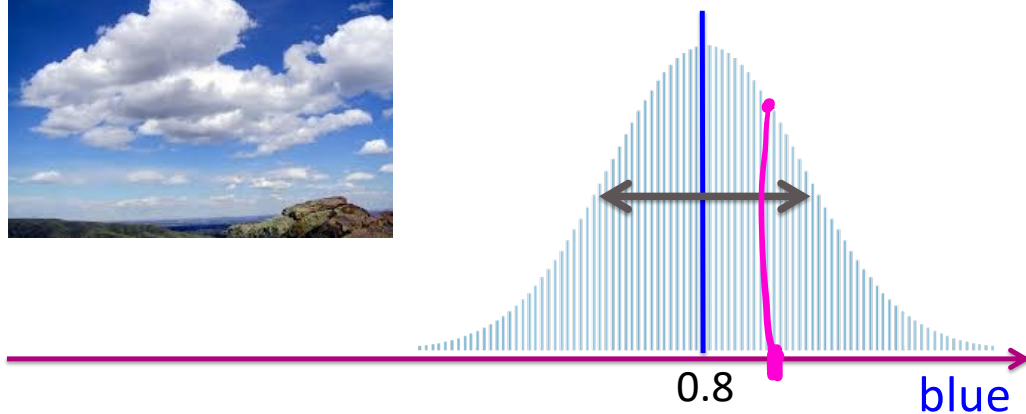
[R = 0.85, G = 0.05, B = 0.35]



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

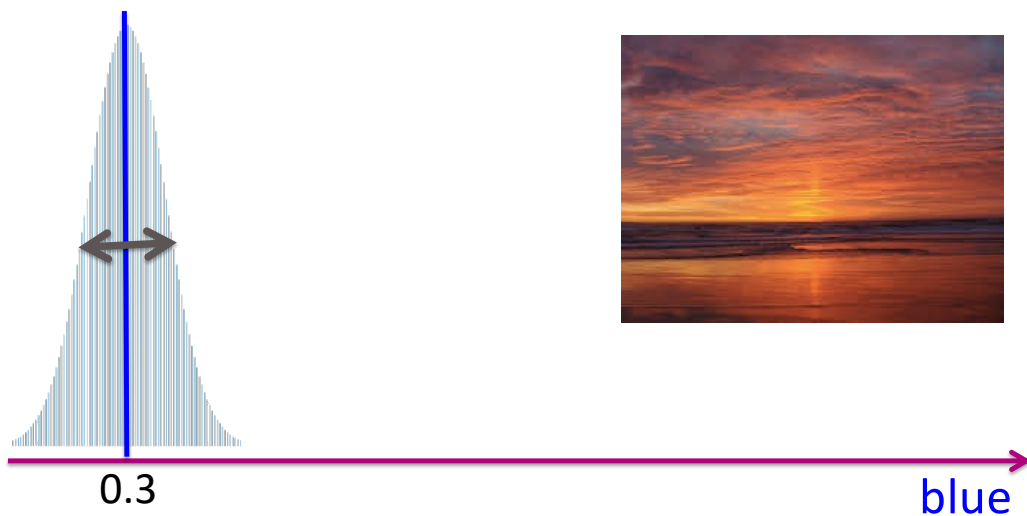
Distribution over all **cloud** images

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



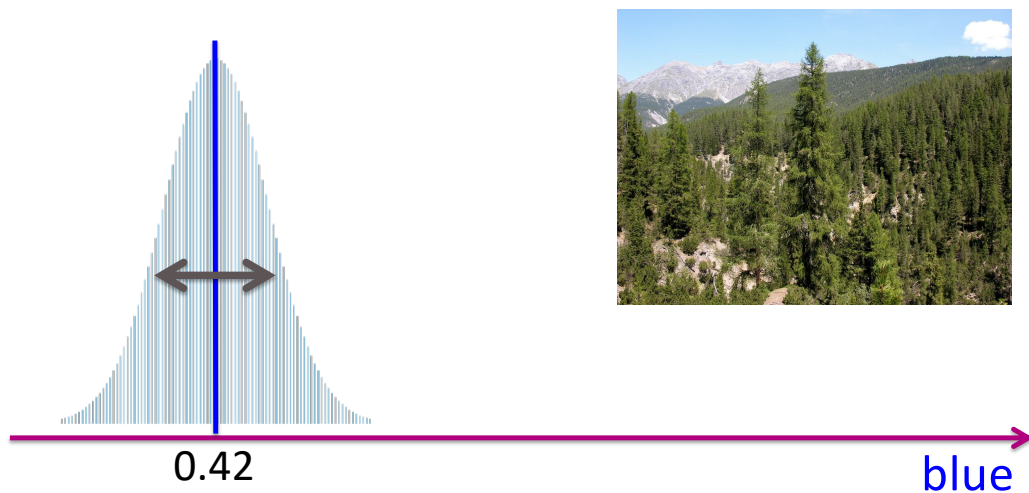
Distribution over all sunset images

Let's look at just the blue dimension

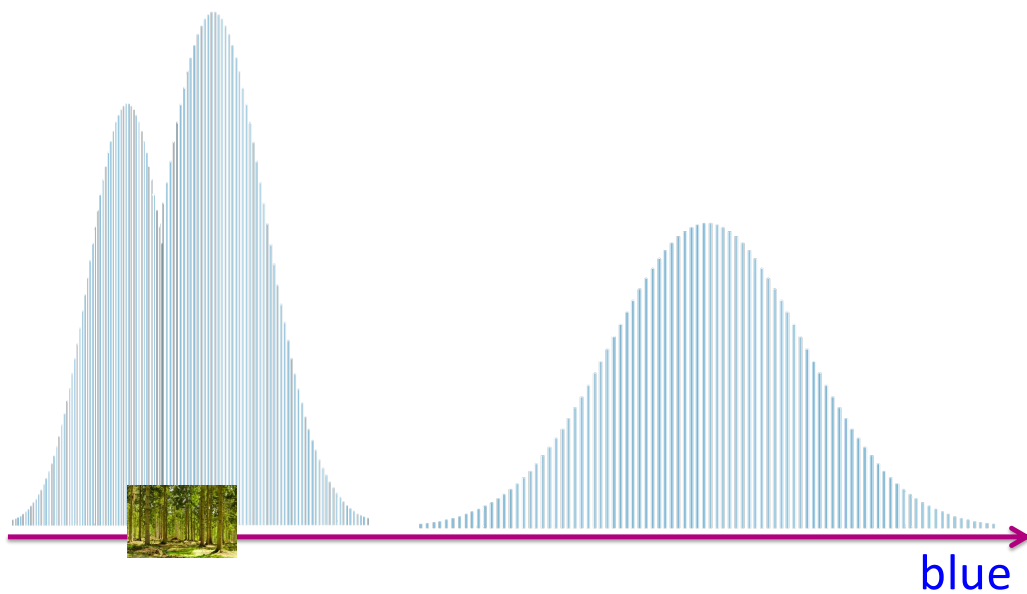


Distribution over all forest images

Let's look at just the blue dimension

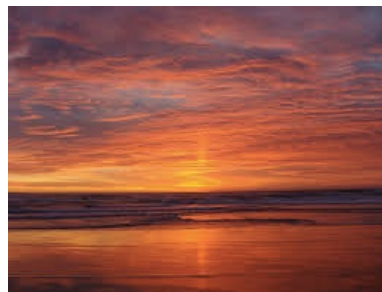
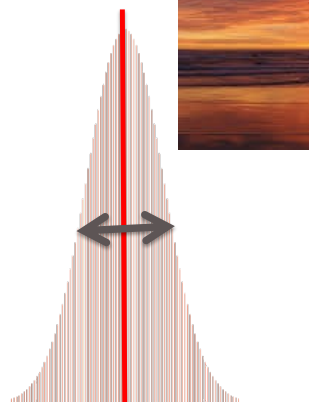
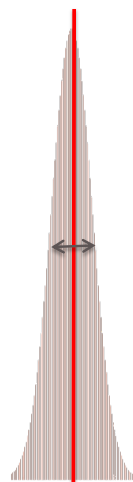


Distribution over **all** images



Can be distinguished along other dim

Now look at the **red** dimension



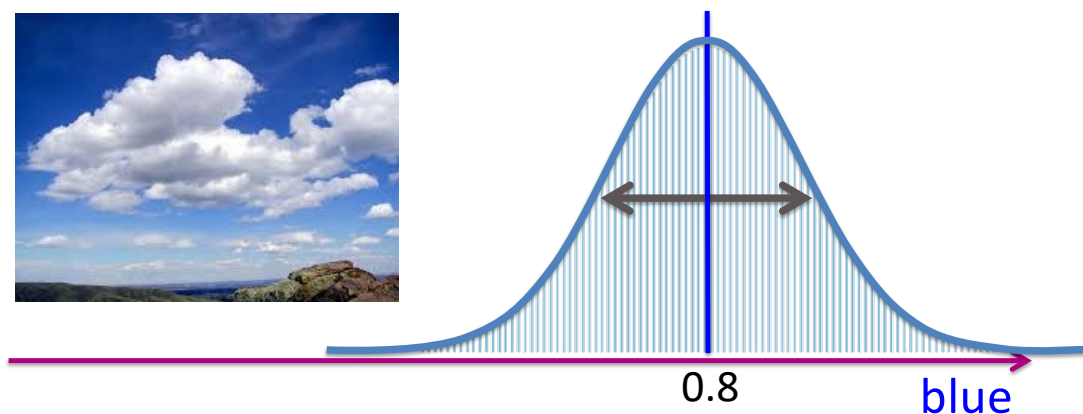
0.05

0.9 **red**

Model for a given image type

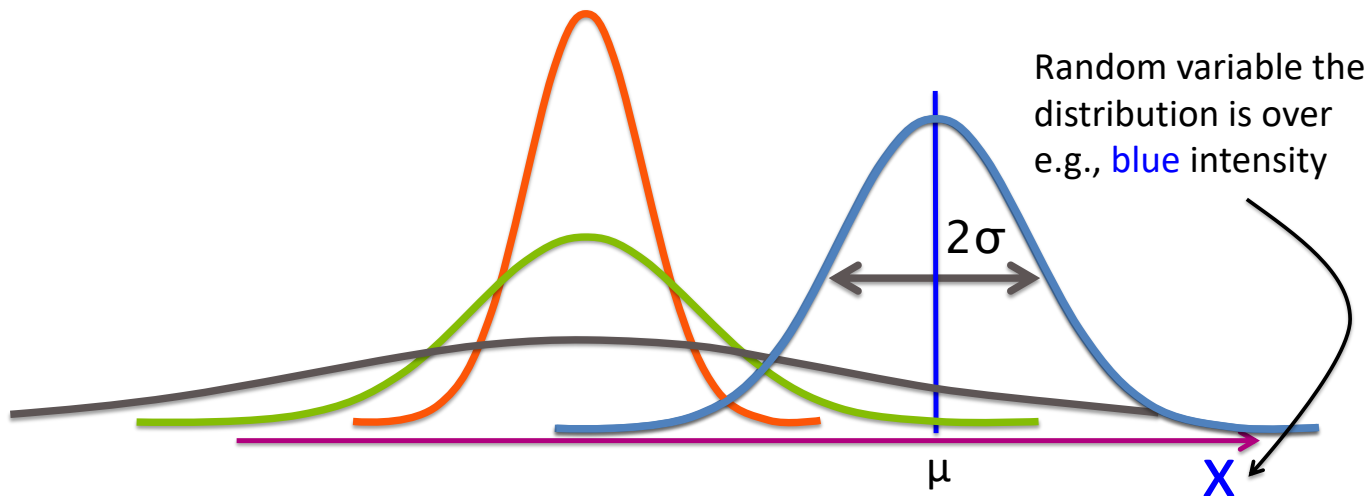


For each dim of the [R, G, B] vector, and each image type, assume a Gaussian distribution over color intensity



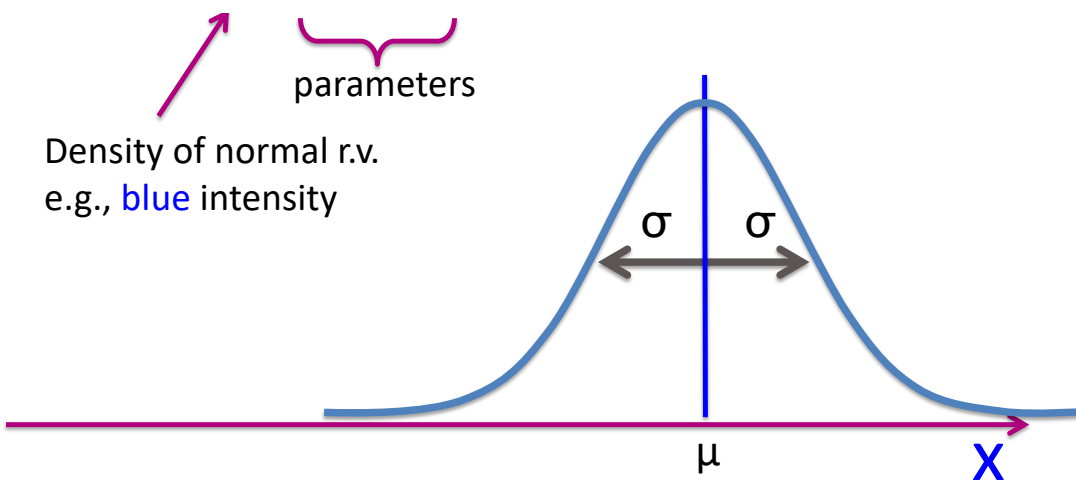
1D Gaussians

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



Density of 1D Gaussian distribution

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



Covariance Matrix $\sum_{(x,y)} (x - E(X))(y - E(Y)) Pr(\underline{X=x, Y=y})$

Problem 1 on your current homework

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

$= \text{Cov}(Y, X)$

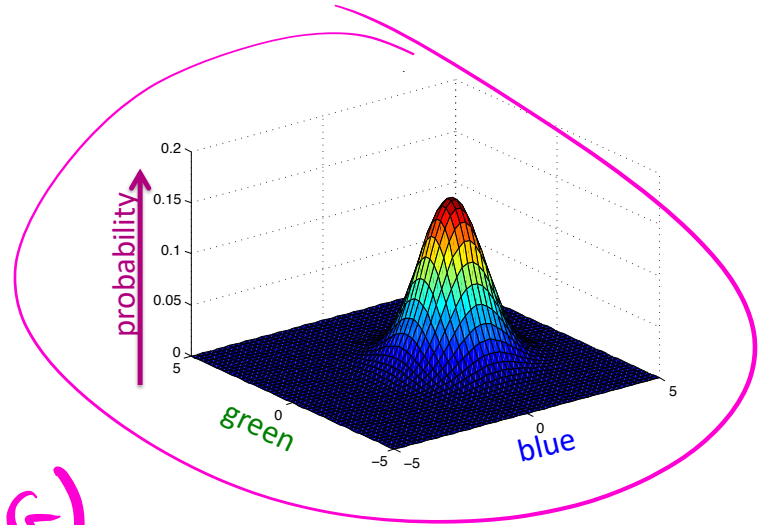
Given set of random variables X_1, X_2, \dots, X_n

The “**covariance matrix**”

$$\Sigma = \begin{matrix} & & j & & \\ \begin{matrix} i \\ \vdots \\ \vdots \end{matrix} & \begin{bmatrix} \text{Cov}(X_1, X_1) & \dots & \text{Cov}(X_1, X_n) \\ \vdots & \boxed{\text{Cov}(X_i, X_i)} & \vdots \\ \text{Cov}(X_n, X_1) & \dots & \text{Cov}(X_n, X_n) \end{bmatrix} \end{matrix}$$

Bivariate Gaussian (2 dimensions)

Fully specified by means μ
and covariance matrix Σ



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

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

$\text{Cov}(B, G)$

$\sigma_{\text{green green}}$

covariance determines
orientation + spread

$\text{Var}(B)$
 $\text{Cov}(B, B)$
 $\text{Cov}(B, G)$
 $\text{Cov}(G, B)$
 $\text{Cov}(G, G)$
 $\text{Var}(G)$

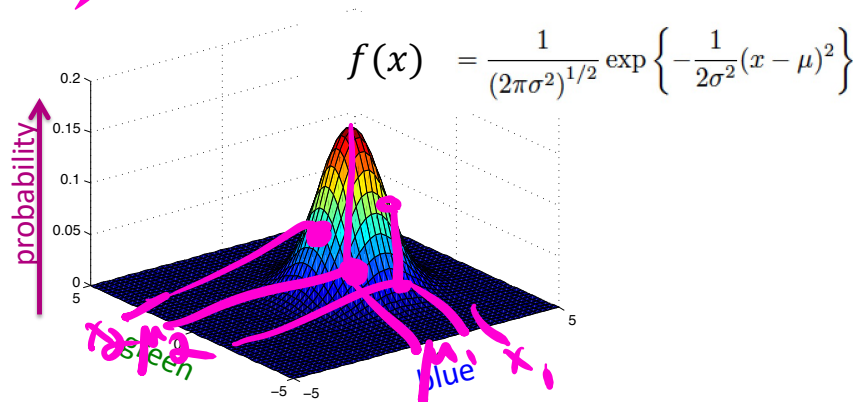
Multivariate Gaussian (example – bivariate)

Fully specified by means μ
and covariance matrix Σ

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

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

covariance determines
orientation + spread



x, y

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

$\exp\{ \} \equiv e^{\{ \}}$

$$-\frac{1}{2} (x-\mu) \left(\frac{1}{\sigma^2} \right) (x-\mu)$$

(x_1, x_2)

Multi-dimensional

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

1-dimensional

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

$$\boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

(x)

$$(x_1 - \mu_1, x_2 - \mu_2) (\boldsymbol{\Sigma}^{-1}) \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}$$

$$|2\pi\boldsymbol{\Sigma}| = \det(2\pi\boldsymbol{\Sigma}) = (2\pi)^n \det(\boldsymbol{\Sigma})$$

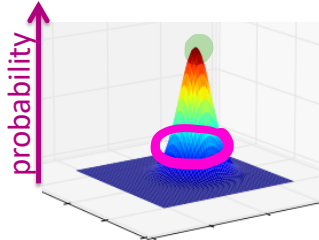
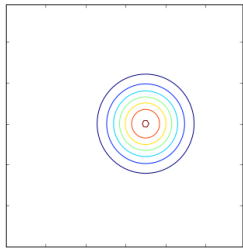
$\boldsymbol{\Sigma}$ is $n \times n$

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

scalar.

Covariance structure

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



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

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

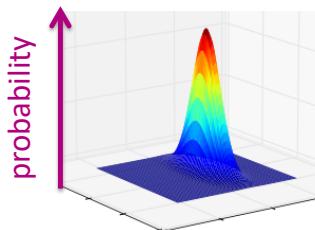
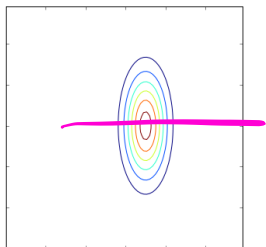
$$\Sigma^{-1} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{\sigma^2} \end{pmatrix}$$

$$\begin{pmatrix} x_1 & x_2 \end{pmatrix} \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{\sigma^2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \frac{x_1^2}{\sigma^2} + \frac{x_2^2}{\sigma^2} = \frac{x_1^2 + x_2^2}{\sigma^2}$$

Covariance structure

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

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

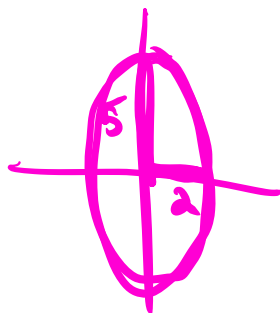


$$\begin{pmatrix} \sigma^2 & 0 \\ 0 & 5\sigma^2 \end{pmatrix}$$

$$\boldsymbol{\Sigma}^{-1} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{5\sigma^2} \end{pmatrix}$$

$$(x_1, x_2) \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{5\sigma^2} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$\frac{x_1^2}{\sigma^2} + \frac{x_2^2}{5\sigma^2} = 1$$



$$\begin{pmatrix} 5 & 0 \\ 0 & 1 \end{pmatrix}$$

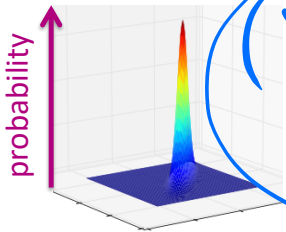
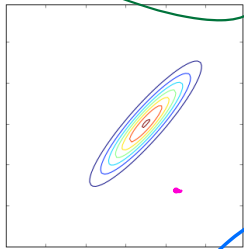
$$Ax = \lambda x$$

(0 3) (1 0)

Covariance Structure

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

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

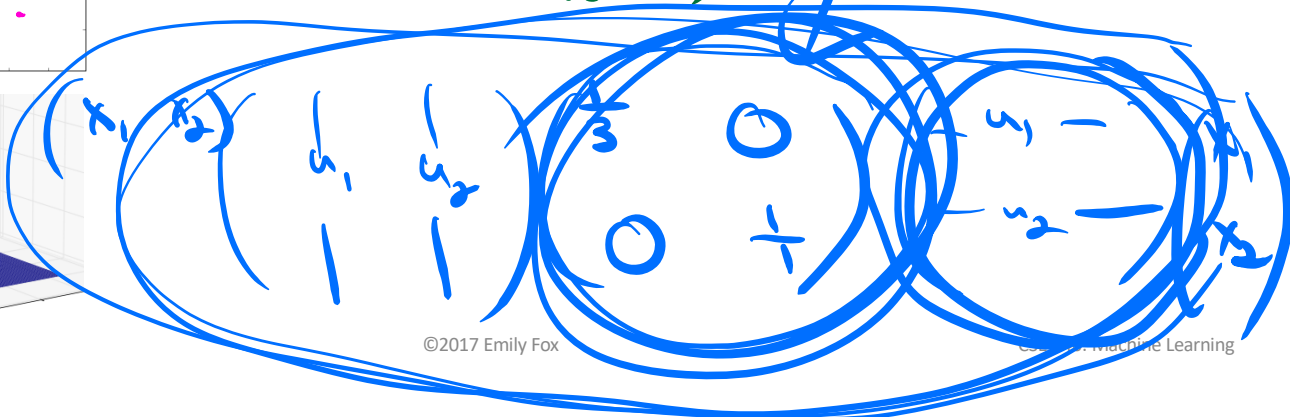


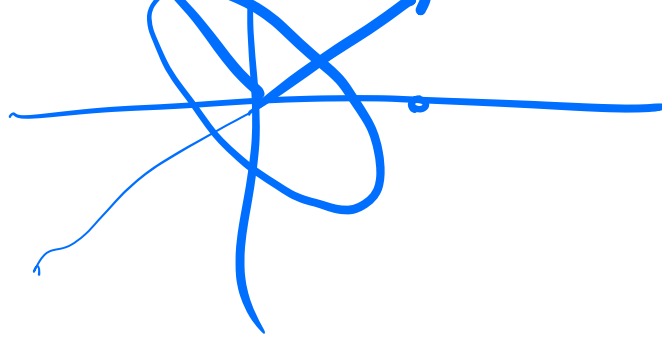
$$\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

$$\lambda_1 = 3 \quad u_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\lambda_2 = 1 \quad u_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

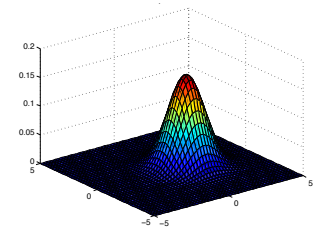
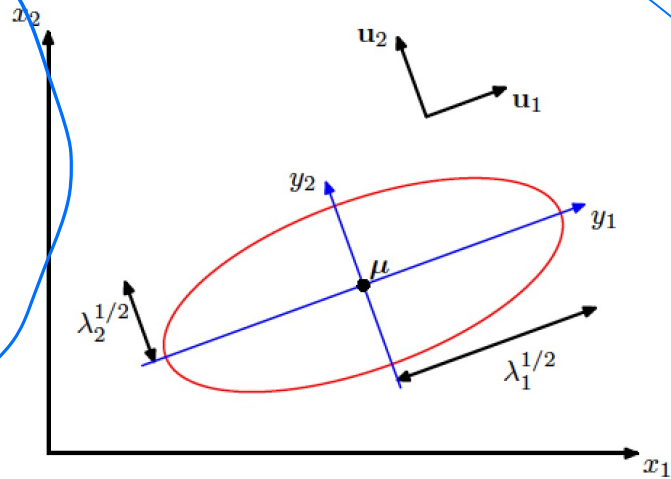
$$\Sigma^{-1} = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$$





Advanced...

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



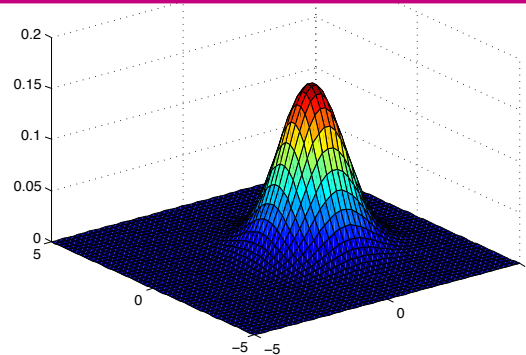
Multivariate Gaussian density

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

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

$$f(\mathbf{x} | \underbrace{\boldsymbol{\mu}, \boldsymbol{\Sigma}}_{\text{parameters}})$$

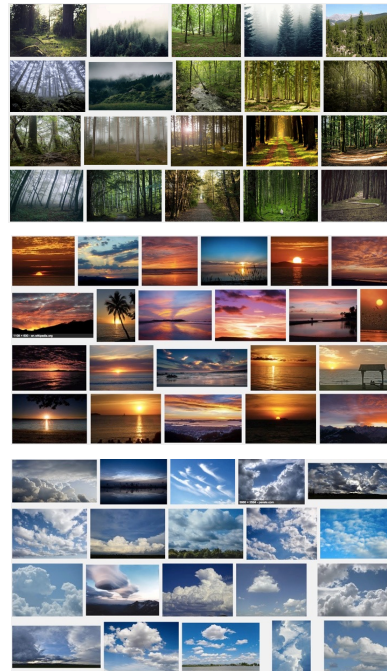
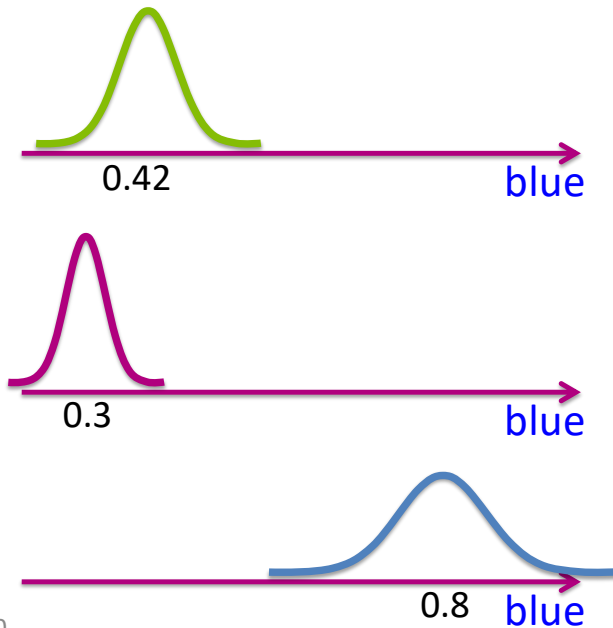
Random vector
e.g., [R, G, B] intensities





Clustering using mixture model

Model as Gaussian per category/cluster

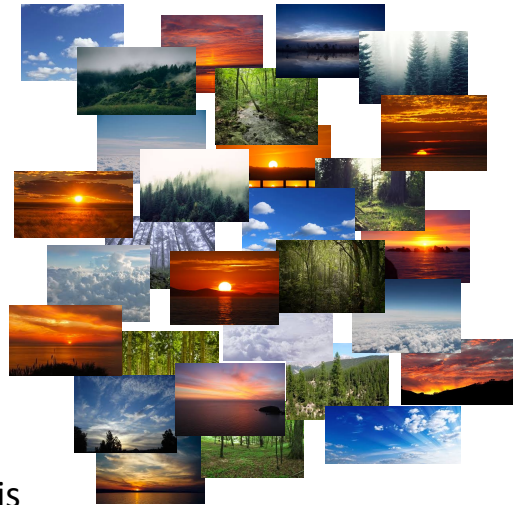
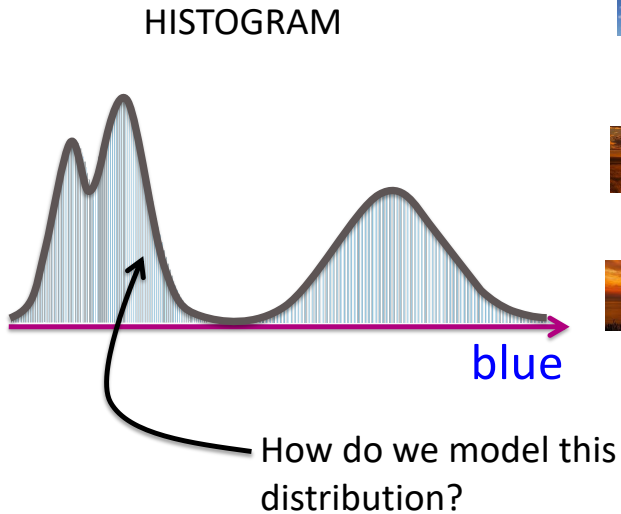


Forests

Sunsets

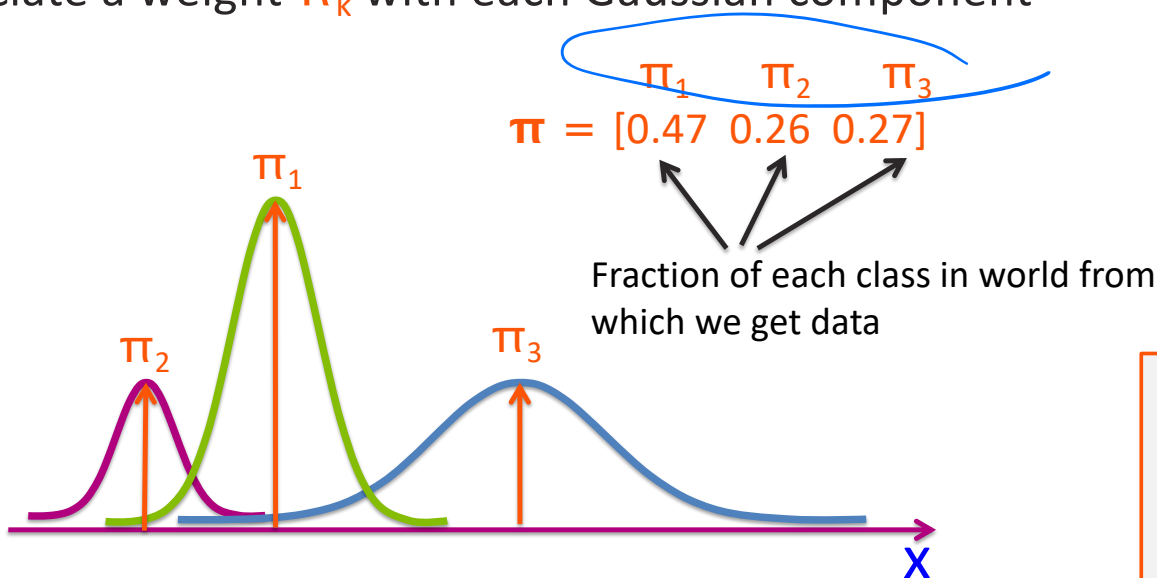
Clouds

Jumble of unlabeled images



Combination of weighted Gaussians

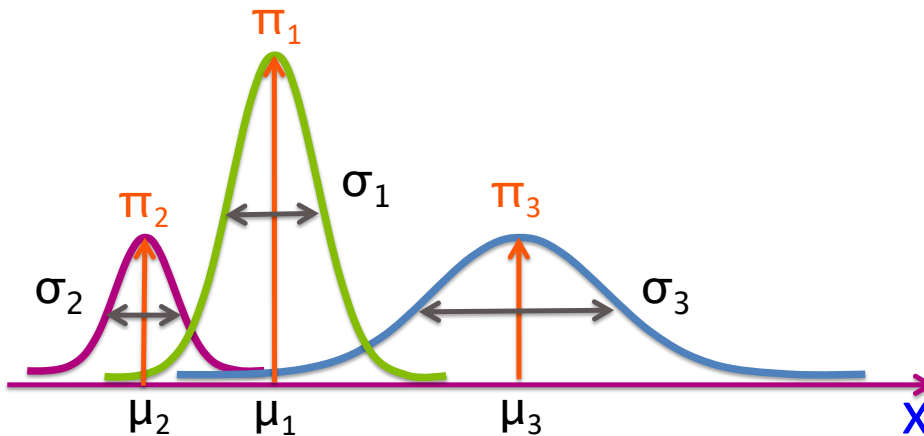
Associate a weight π_k with each Gaussian component



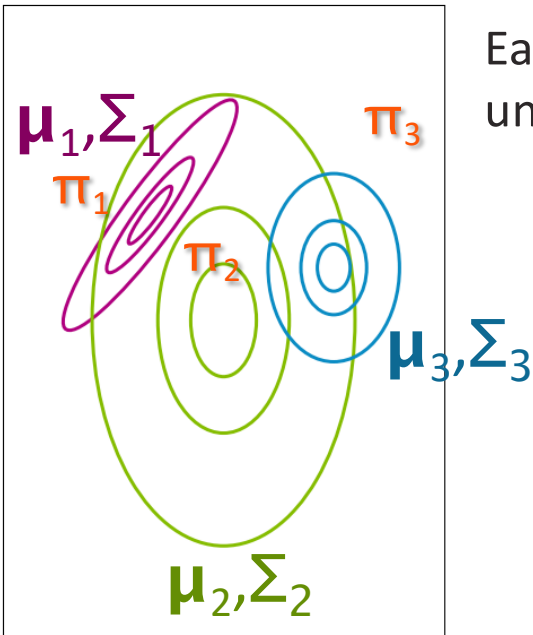
$$0 \leq \pi_k \leq 1$$
$$\sum_{k=1}^K \pi_k = 1$$

Mixture of Gaussians (1D)

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



Mixture of Gaussians (general)



Each mixture component represents a unique cluster specified by:

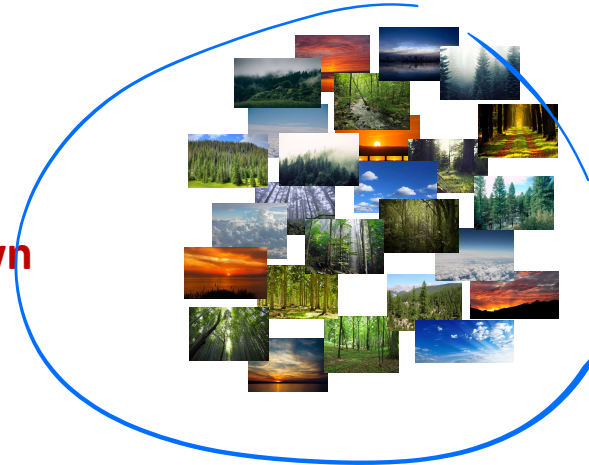
$$\{\pi_k, \mu_k, \Sigma_k\}$$

Mixture model

- K clusters, defined by the following **unknown** parameters

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

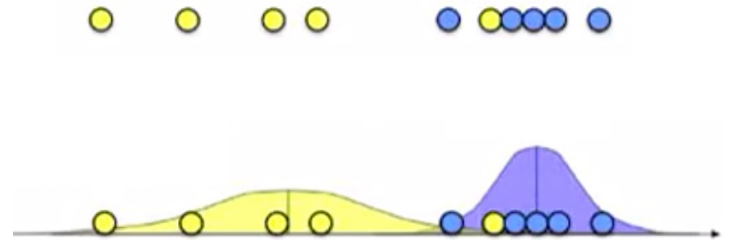
$$\sum_{j=1}^k \pi_j = 1.$$



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

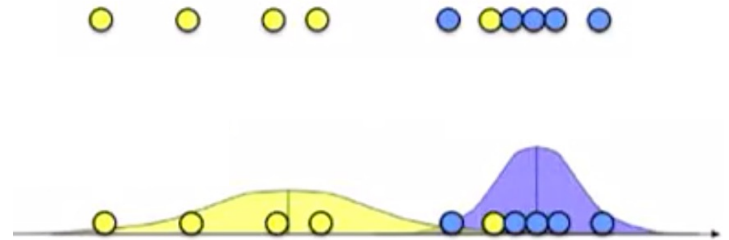
Two 1-D Gaussians, with unknown mean and variance

- Easy if know the source of each data point.



Two 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?



Mixture model

- K clusters, defined by the following parameters

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

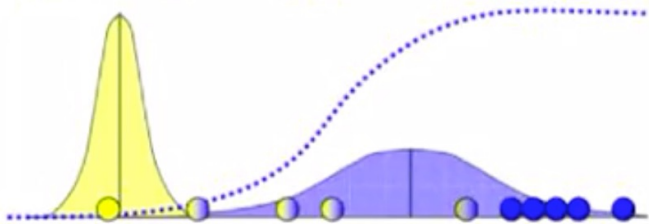
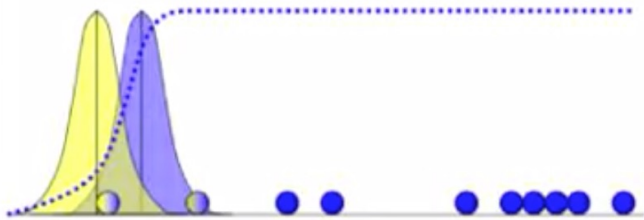
- Problem: Assume that the data comes from such a distribution, and recover the parameters of the distribution.
- Determine, for each point, the likelihood of it belonging to cluster j, for each j.
- **PROBLEM: no closed form solution**



Expectation Maximization Algorithm

Two step approach based on following observation

- If we knew which cluster each sample was from, we could estimate all the parameters.
- If we knew all the parameters we could estimate the chance each point came from each cluster.
- EM is an iterative algorithm that alternates between these two steps.



©2017 Emily Fox

CSE 446: Machine Learning

