

# EXPECTATION MAXIMIZATION

J. Elder

CSE 6390/PSYC 6225 Computational Modeling of Visual Perception

# Credits

- Some of these slides were sourced and/or modified from:
  - ▣ Christopher Bishop, Microsoft UK
  - ▣ Simon Prince, University College London

# Mixtures of Gaussians

3

Expectation Maximization

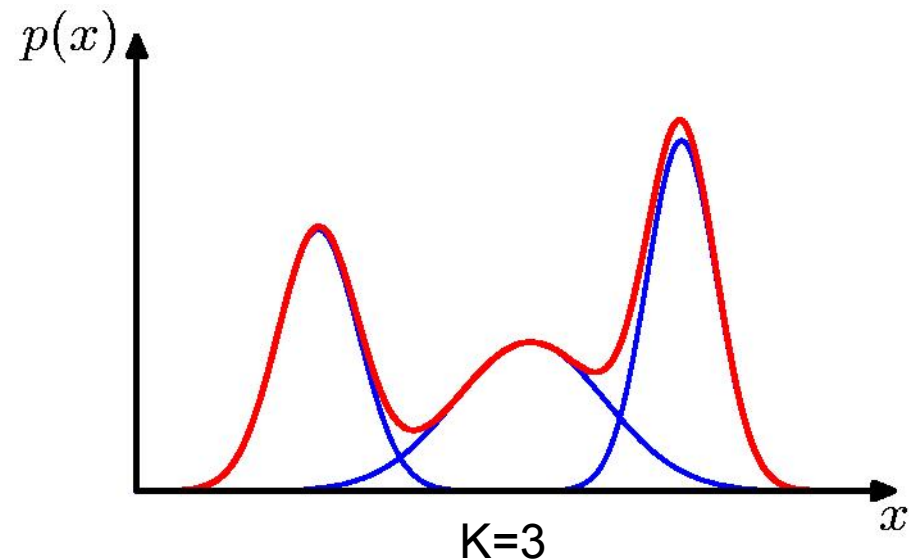
- Combine simple models into a complex model:

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

↑  
Mixing coefficient

Component

$$\forall k : \pi_k \geq 0 \quad \sum_{k=1}^K \pi_k = 1$$



# Mixtures of Gaussians

- Determining parameters  $\mu$ ,  $\sigma$  and  $\pi$  using maximum log likelihood

$$\ln p(\mathbf{X}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left\{ \underbrace{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{\text{Log of a sum; no closed form maximum.}} \right\}$$

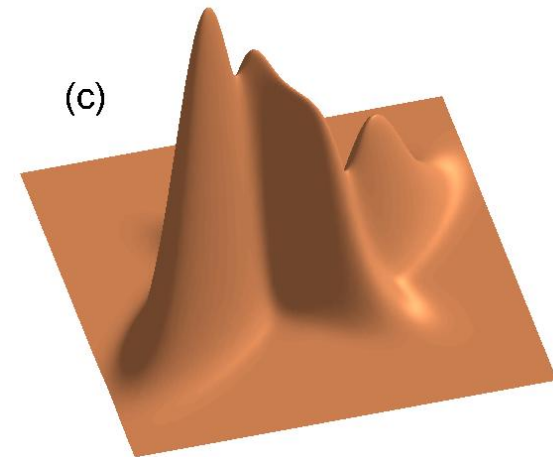
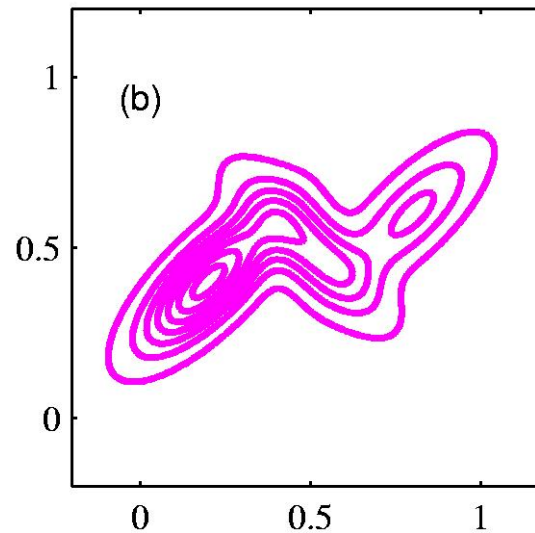
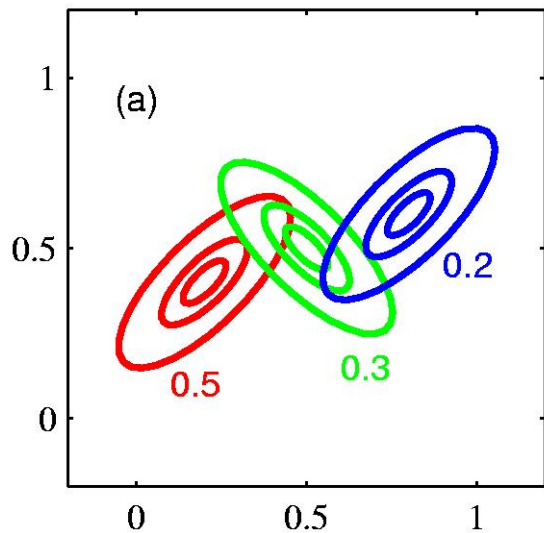
Log of a sum; no closed form maximum.

- Solution: use standard, iterative, numeric optimization methods or the *expectation maximization* algorithm (Chapter 9).

# Mixtures of Gaussians

5

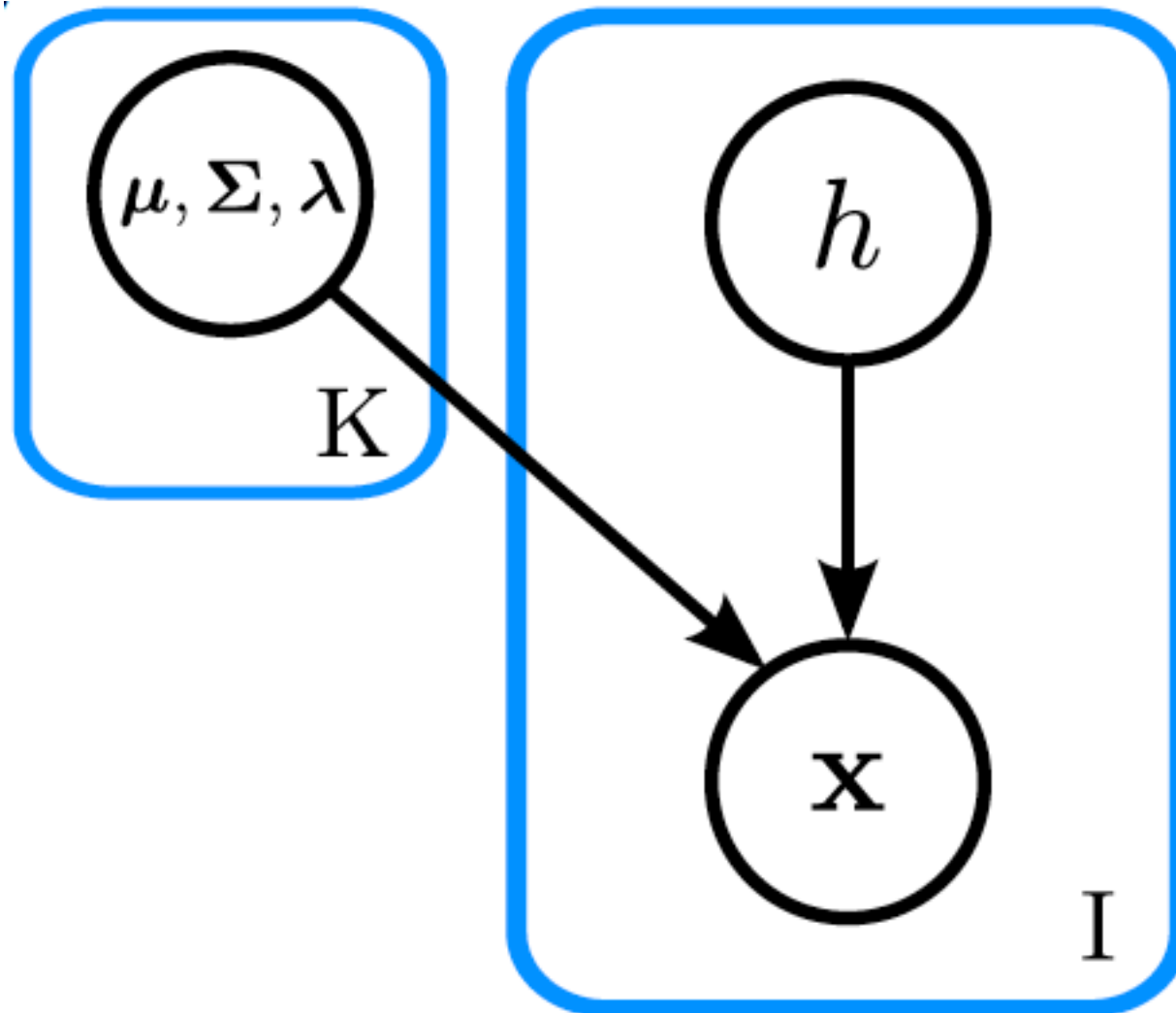
Expectation Maximization



# Graphical Model for Gaussian Mixture

6

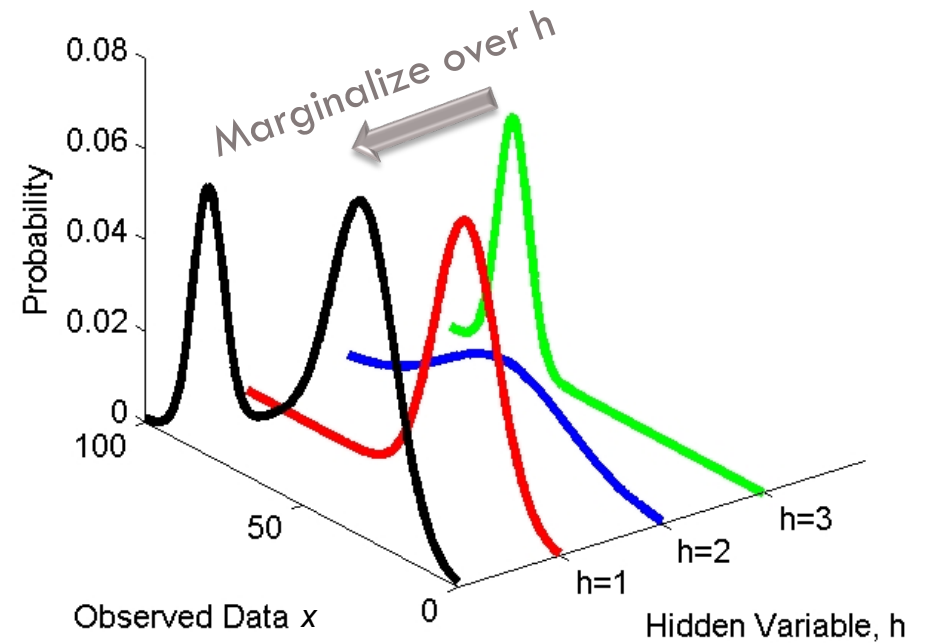
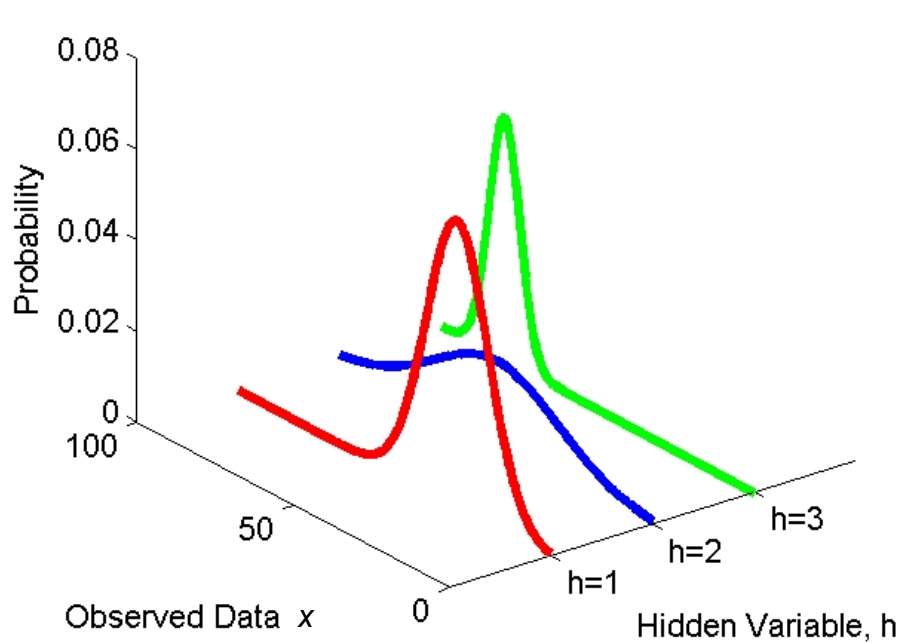
Expectation Maximization



# Hidden Variable Interpretation

$$Pr(x|w_{1...K}, \mu_{1...K}, \sigma_{1...K}^2) = \sum_{k=1}^K w_k \mathcal{G}_x [\mu_k, \sigma_k^2]$$

$$= \sum_{k=1}^K Pr(h = k) Pr(x|h = k)$$



# Hidden Variable Interpretation

8

Expectation Maximization

## ASSUMPTIONS

- for each training datum  $x_i$  there is a hidden variable  $h_i$ .
- $h_i$  represents which Gaussian  $x_i$  came from
- hence  $h_i$  takes discrete values

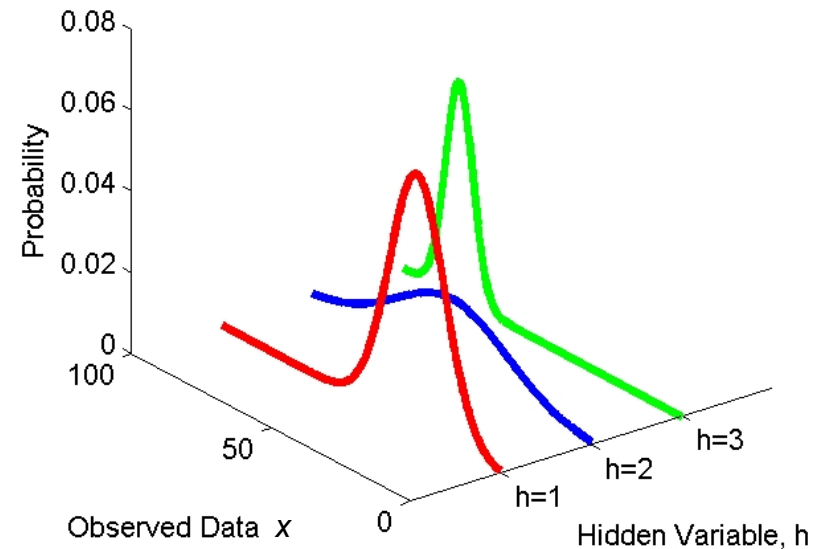
## OUR GOAL:

To estimate the parameters  $\theta$ :

The means  $\mu_k$

The covariances  $\Sigma_k$

The weights (mixing coefficients)  $w_k$   
for all  $K$  components of the model.



## THING TO NOTICE:

If we knew the hidden variables  $h_i$  for the training data it would very easy to estimate parameters  $\theta$  – just estimate individual Gaussians separately.



# Hidden Variable Interpretation

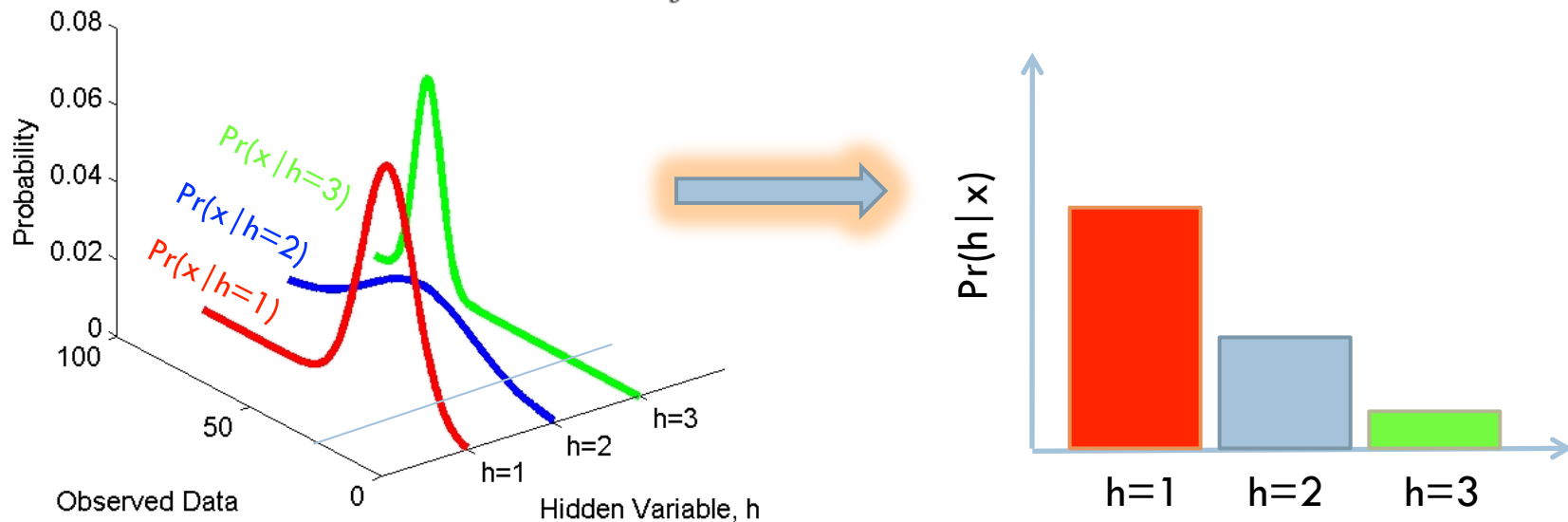
9

Expectation Maximization

THING TO NOTICE #2:

If we knew the parameters  $\theta$  it would be very easy to estimate the posterior distribution over each hidden variable  $h_i$  using Bayes' rule:

$$Pr(h_i = k | x_i, \theta) = \frac{Pr(x_i | h_i = k, \theta) Pr(h_i = k)}{\sum_{j=1}^K Pr(x_i | h_i = j, \theta) Pr(h_i = j)}$$



# Expectation Maximization

10

Expectation Maximization

## Chicken and egg problem:

- could find  $h_{1\dots N}$  if we knew  $\theta$
- could find  $q$  if we knew  $h_{1\dots N}$

## Solution: Expectation Maximization (EM) algorithm

(Dempster, Laird and Rubin 1977)

Alternate between:

### 1. Expectation Step (E-Step)

- For fixed  $\theta$  find posterior distribution over  $h_{1\dots N}$

### 2. Maximization Step (M-Step)

- Given these distributions, maximize lower bound on likelihood w.r.t.  $\theta$

# Expectation Maximization

$$Pr(\mathbf{x}|\boldsymbol{\theta}) = \int Pr(\mathbf{x}, \mathbf{h}|\boldsymbol{\theta})d\mathbf{h} = \int Pr(\mathbf{x}|\mathbf{h}, \boldsymbol{\theta})Pr(\mathbf{h})d\mathbf{h}$$

We introduce a probability distribution  $q(\mathbf{h})$  over the hidden variables  $\mathbf{h}$ .

EM works by defining a lower bound  $B[q(\mathbf{h}), \boldsymbol{\theta}]$  on the log likelihood  $\log P(\mathbf{x} | \boldsymbol{\theta})$ .

and then iteratively increasing this lower bound by alternately updating

$q(\mathbf{h})$  (E-step)

and

$\boldsymbol{\theta}$  (M-step)

# E-Step

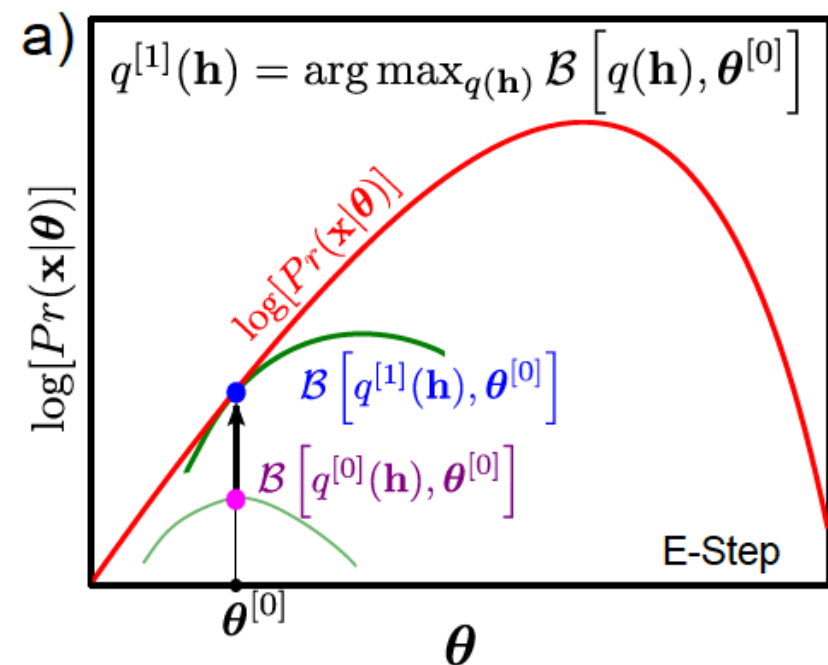
12

Expectation Maximization

Fix  $\theta$ .

Find  $q(\mathbf{h})$  that maximizes lower bound:

$$q_i^{[t]}[\mathbf{h}] = \arg \max_{q_i[\mathbf{h}]} \mathcal{B}[q_i[\mathbf{h}], \theta^{[t-1]}].$$



# M-Step

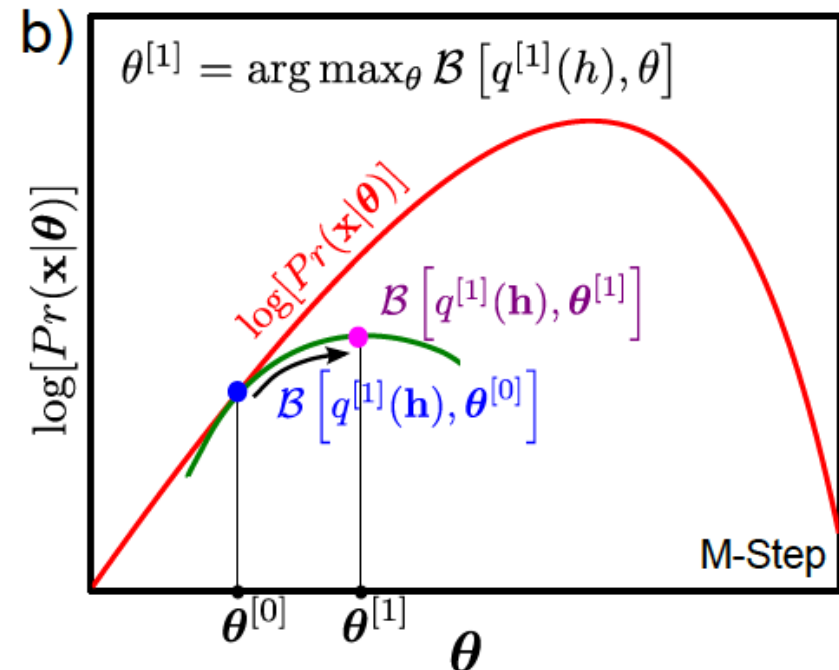
13

Expectation Maximization

Fix  $q(\mathbf{h})$ .

Find  $\theta$  that maximizes lower bound:

$$\theta^{[t]} = \arg \max_{\theta} \mathcal{B}[q^{[t]}(\mathbf{h}), \theta].$$



# Lower Bound on Likelihood

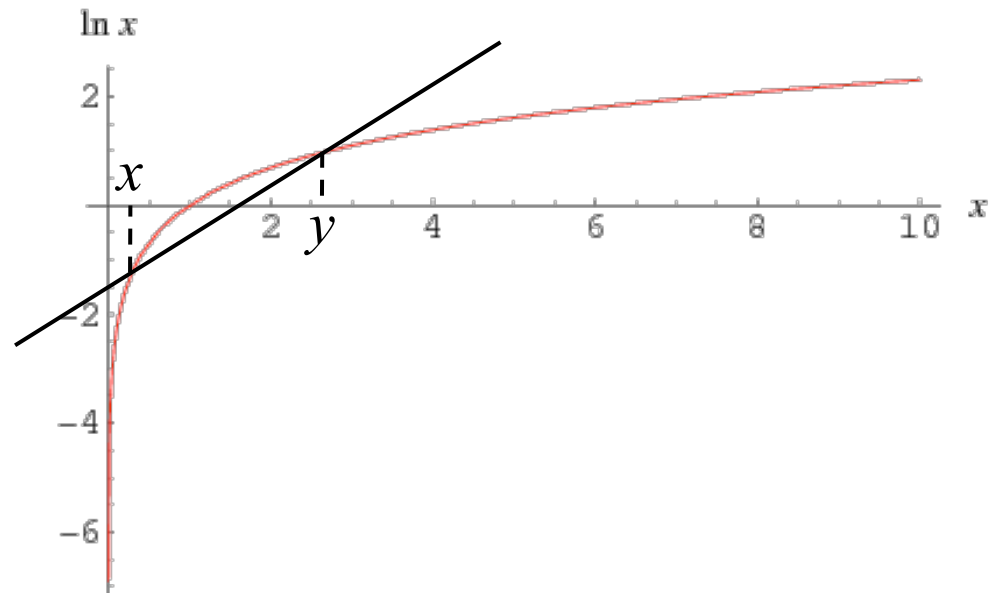
$$\mathcal{B}[q_i(\mathbf{h}_i), \theta] = \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{Pr(\mathbf{x}, \mathbf{h}_i | \theta)}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_{1\dots I}$$

Note that  $\log$  is a concave function, i.e.,

$$\log(tx + (1-t)y) \geq tf(x) + (1-t)f(y) \quad \forall x, y > 0, 0 \leq t \leq 1$$

or equivalently

$$\frac{\partial^2}{\partial x^2} \log x \leq 0$$



# Jensen's Inequality

15

Expectation Maximization

If  $f(x)$  is a concave function, then  $E[f(x)] \leq f(E[x])$

Proof: Note that  $f(y) - f(x) \geq (y - x)f'(y) \quad \forall y > 0$

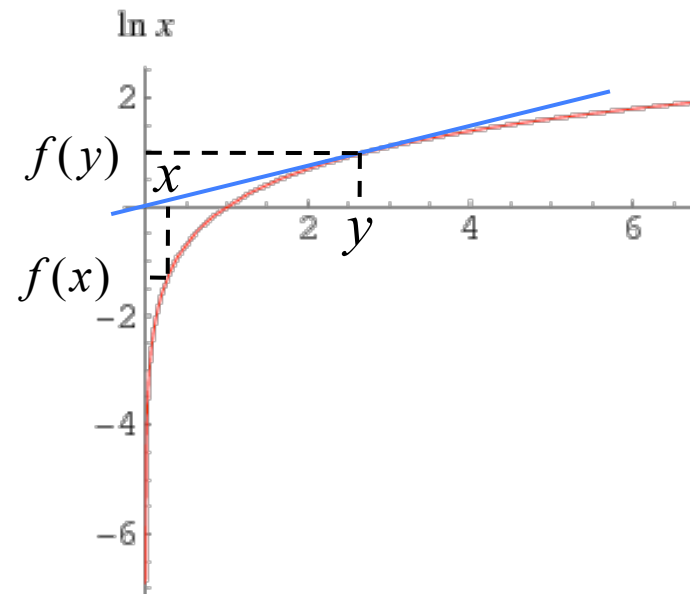
Choose  $y = E[x]$ . Then

$$f(E[x]) - f(x) \geq (E[x] - x)f'(E[x])$$

Now, taking expectations of both sides,

$$f(E[x]) - E[f(x)] \geq 0$$

$$\rightarrow E[f(x)] \leq f(E[x])$$



# Lower Bound on Likelihood

Thus

$$\begin{aligned}\mathcal{B}[q_i(\mathbf{h}_i), \boldsymbol{\theta}] &= \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{Pr(\mathbf{x}, \mathbf{h}_i | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_{1\dots I} \\ &\leq \sum_{i=1}^I \log \left[ \int q_i(\mathbf{h}_i) \frac{Pr(\mathbf{x}, \mathbf{h} | \boldsymbol{\theta})}{q_i(\mathbf{h}_i)} d\mathbf{h} \right] \\ &= \sum_{i=1}^I \log \left[ \int Pr(\mathbf{x}, \mathbf{h} | \boldsymbol{\theta}) d\mathbf{h} \right] . \\ &= \log P(\mathbf{x} | \boldsymbol{\theta})\end{aligned}$$



# E-Step

17

Expectation Maximization

□ How do we maximize the bound in the E-step?

$$\begin{aligned}\mathcal{B}[q_i(\mathbf{h}_i), \theta] &= \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{Pr(\mathbf{x}, \mathbf{h}_i | \theta)}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_{1\dots I} \\ &= \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{Pr(\mathbf{h}_i | \mathbf{x}_i, \theta) Pr(\mathbf{x}_i, \theta)}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_{1\dots I} \\ &= \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log [Pr(\mathbf{x}_i, \theta)] d\mathbf{h}_{1\dots I} - \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{q_i(\mathbf{h}_i)}{Pr(\mathbf{h}_i | \mathbf{x}_i, \theta)} \right] d\mathbf{h}_{1\dots I} \\ &= \underbrace{\sum_{i=1}^I \log [Pr(\mathbf{x}_i, \theta)]}_{\text{Log-Likelihood}} - \underbrace{\sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{q_i(\mathbf{h}_i)}{Pr(\mathbf{h}_i | \mathbf{x}_i, \theta)} \right] d\mathbf{h}_{1\dots I}}_{\text{Kullback-Leibler Divergence}}\end{aligned}$$

# Kullback-Leibler Divergence

$$D_{KL}(P \parallel Q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

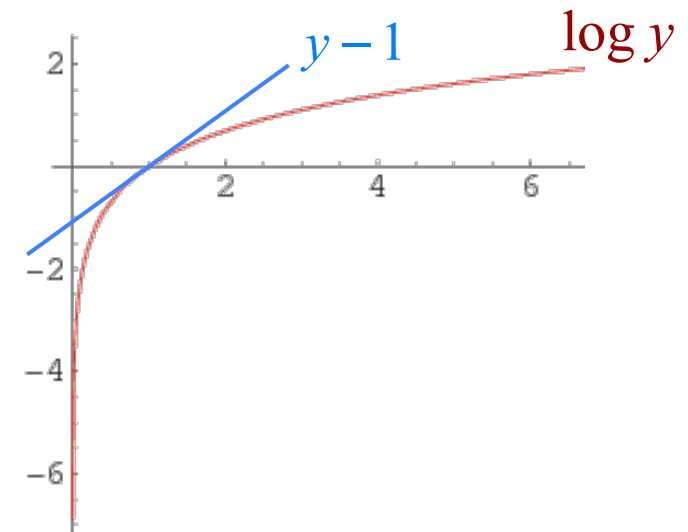
- An information-theoretic measure of the deviation between two distributions
- The KL divergence is strictly non-negative.
  - Proof:

$$D_{KL}(P \parallel Q) = \int p(x) \log \frac{p(x)}{q(x)} dx = - \int p(x) \log \frac{q(x)}{p(x)} dx$$

Note that  $\log y \leq y - 1 \quad \forall y > 0$

$$\text{Thus } \int p(x) \log \frac{q(x)}{p(x)} dx \leq \int p(x) \left( \frac{q(x)}{p(x)} - 1 \right) dx = \int (q(x) - p(x)) dx = 0$$

$$\text{Thus } D_{KL}(P \parallel Q) = \int p(x) \log \frac{p(x)}{q(x)} dx \geq 0.$$



# E-Step

- Thus the bound is maximized when the KL-D is 0.

i.e., when  $q(\mathbf{h}) = P(\mathbf{h} | \mathbf{x}, \theta)$

“Responsibility”:  $h_i$  is responsible for explaining  $x_i$ .

$$\begin{aligned} \mathcal{B}[q_i(\mathbf{h}_i), \theta] &= \sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{Pr(\mathbf{x}, \mathbf{h}_i | \theta)}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_{1\dots I} \\ &= \underbrace{\sum_{i=1}^I \log [Pr(\mathbf{x}_i, \theta)]}_{\text{Log-Likelihood}} - \underbrace{\sum_{i=1}^I \int q_i(\mathbf{h}_i) \log \left[ \frac{q_i(\mathbf{h}_i)}{Pr(\mathbf{h}_i | \mathbf{x}_i, \theta)} \right] d\mathbf{h}_{1\dots I}}_{\text{Kullback-Leibler Divergence}} \end{aligned}$$

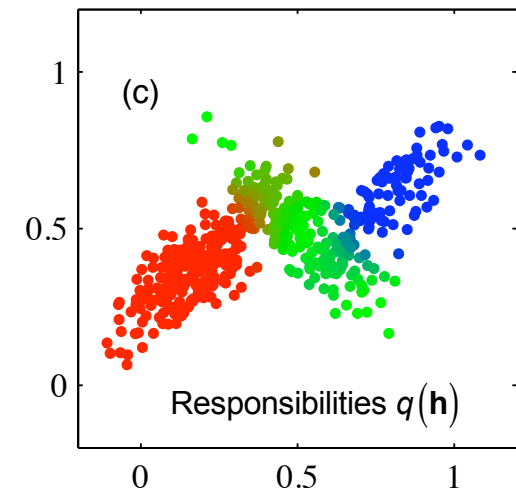
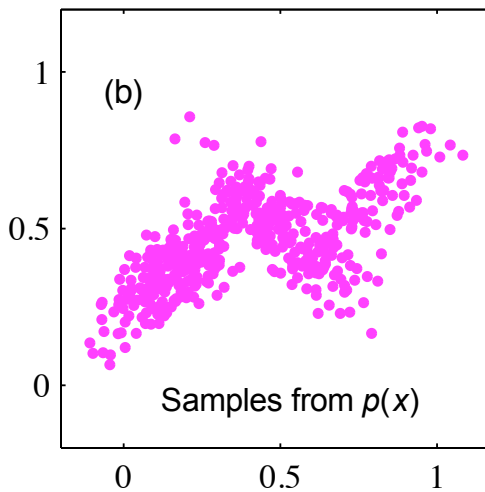
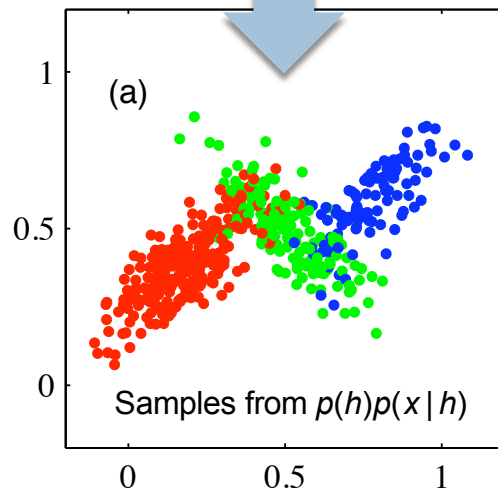
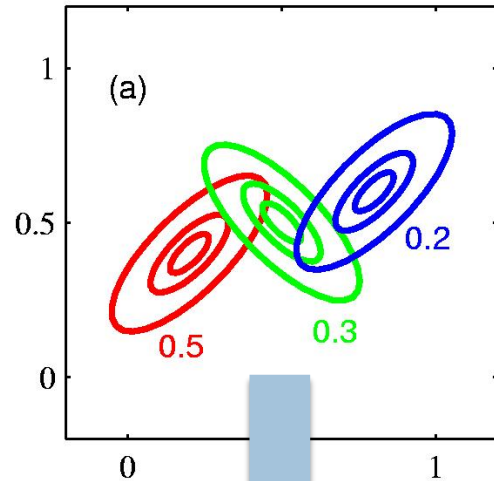
# M-Step

$$\begin{aligned}\theta^{[t]} &= \arg \max_{\theta} \mathcal{B}[q_i^{[t]}(\mathbf{h}_i), \theta] \\ &= \arg \max_{\theta} \sum_{i=1}^I \int q_i^{[t]}(\mathbf{h}_i) \log \left[ \frac{Pr(\mathbf{x}, \mathbf{h}_i | \theta)}{q_i(\mathbf{h}_i)} \right] d\mathbf{h}_{1\dots I} \\ &= \arg \max_{\theta} \sum_{i=1}^I \int q_i^{[t]}(\mathbf{h}_i) \log [Pr(\mathbf{x}, \mathbf{h}_i | \theta)] - q_i(\mathbf{h}_i) \log [q_i(\mathbf{h}_i)] d\mathbf{h}_{1\dots I} \\ &= \arg \max_{\theta} \sum_{i=1}^I \int q_i^{[t]}(\mathbf{h}_i) \log [Pr(\mathbf{x}, \mathbf{h}_i | \theta)] d\mathbf{h}_{1\dots I}.\end{aligned}$$

# Gaussian Mixtures

21

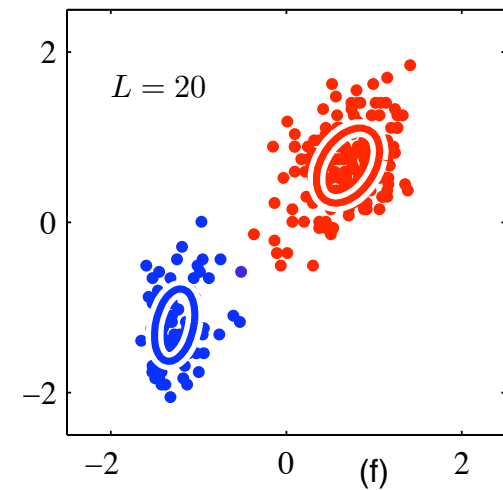
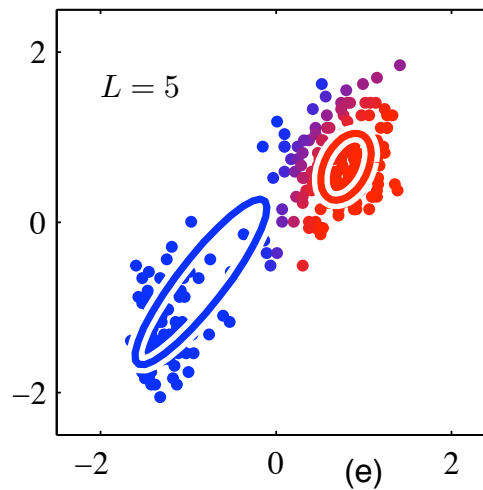
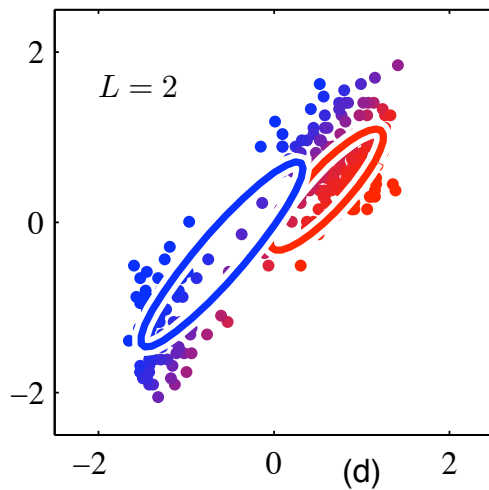
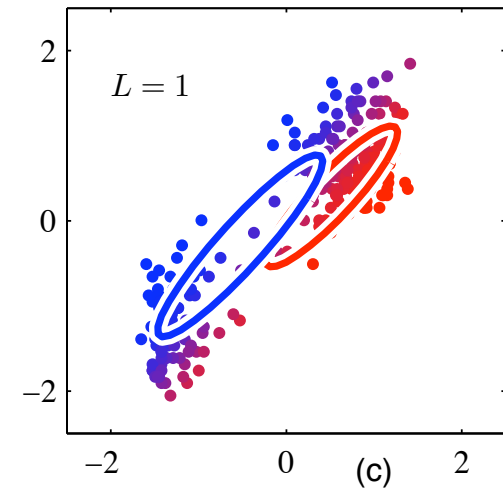
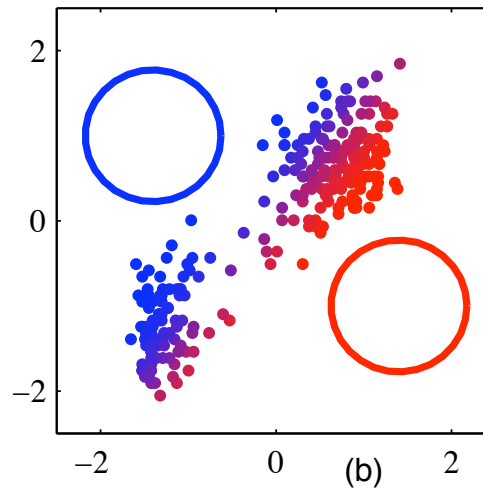
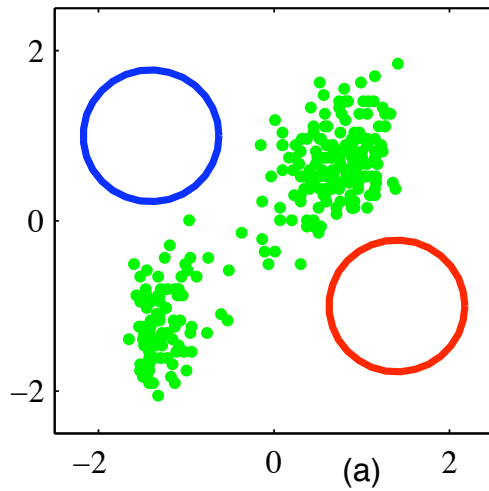
Expectation Maximization



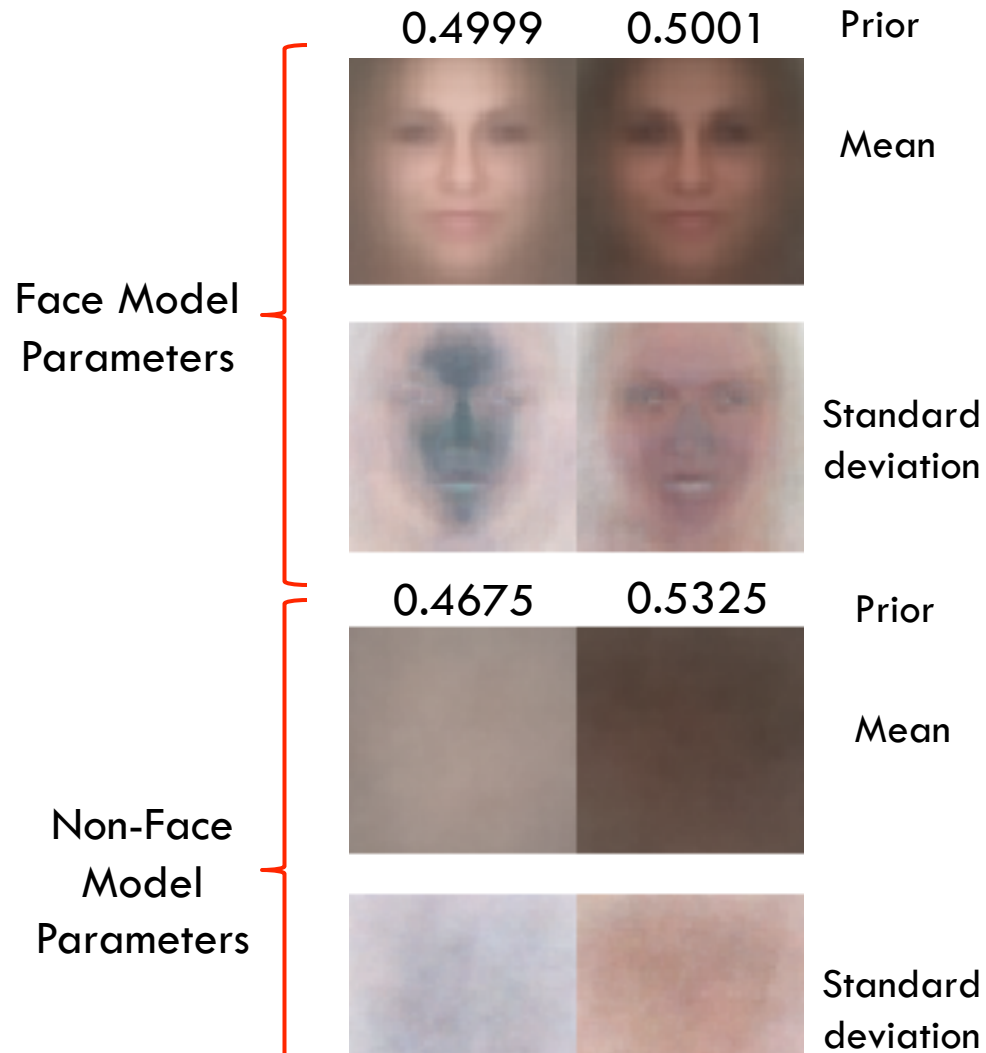
# Old Faithful Example

22

Expectation Maximization



# Face Detection Example: 2 Components

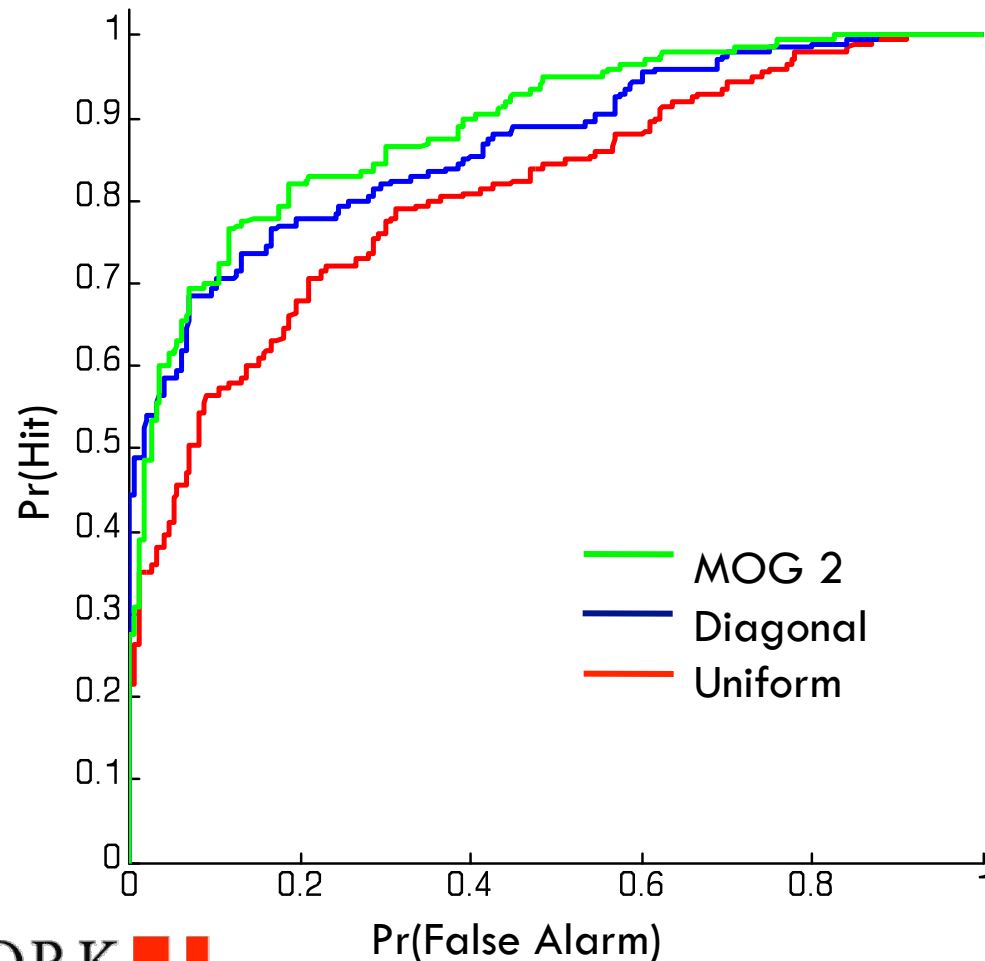


The face model and non-face model have divided the data into two clusters. In each case, these clusters have roughly equal weights.

The primary thing that these seem to have captured is the photometric (luminance) variation.

Note that the standard deviations have become smaller than for the single Gaussian model as any given data point is likely to be close to one mean or the other.

# Results for MOG 2 Model



Performance improves relative to a single Gaussian model, although it is not dramatic.

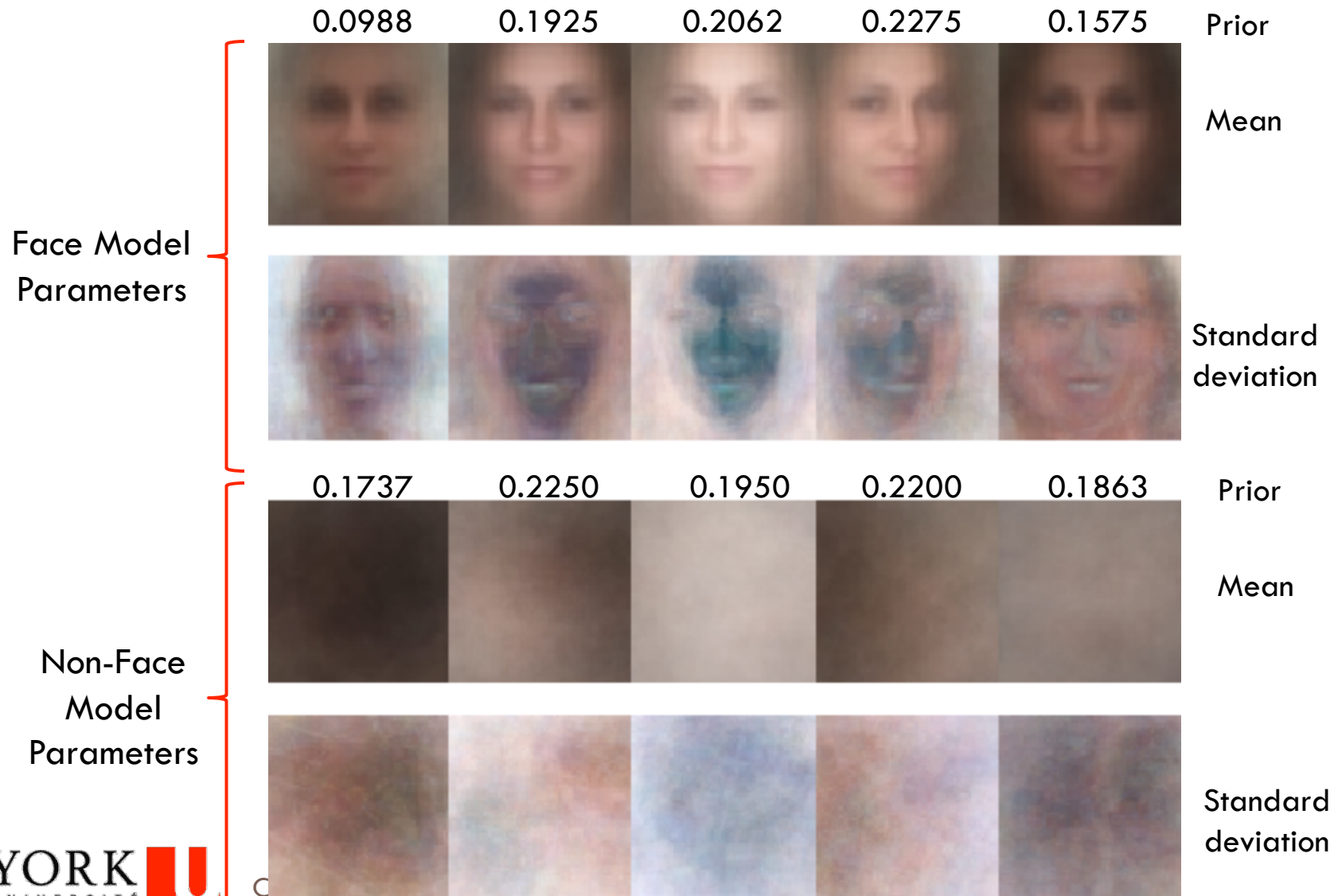
We have a better description of the data likelihood.



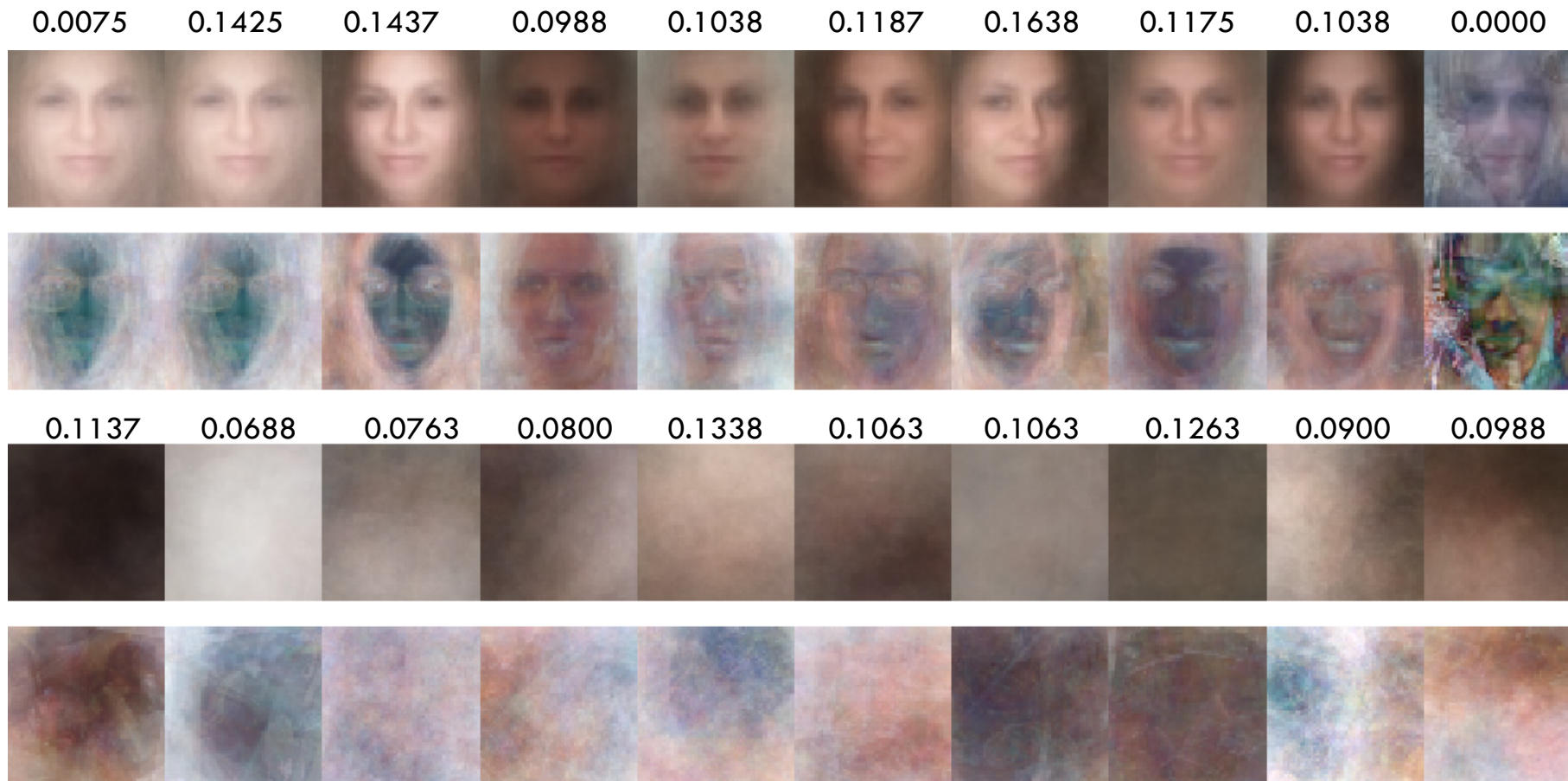
# MOG 5 Components

25

Expectation Maximization



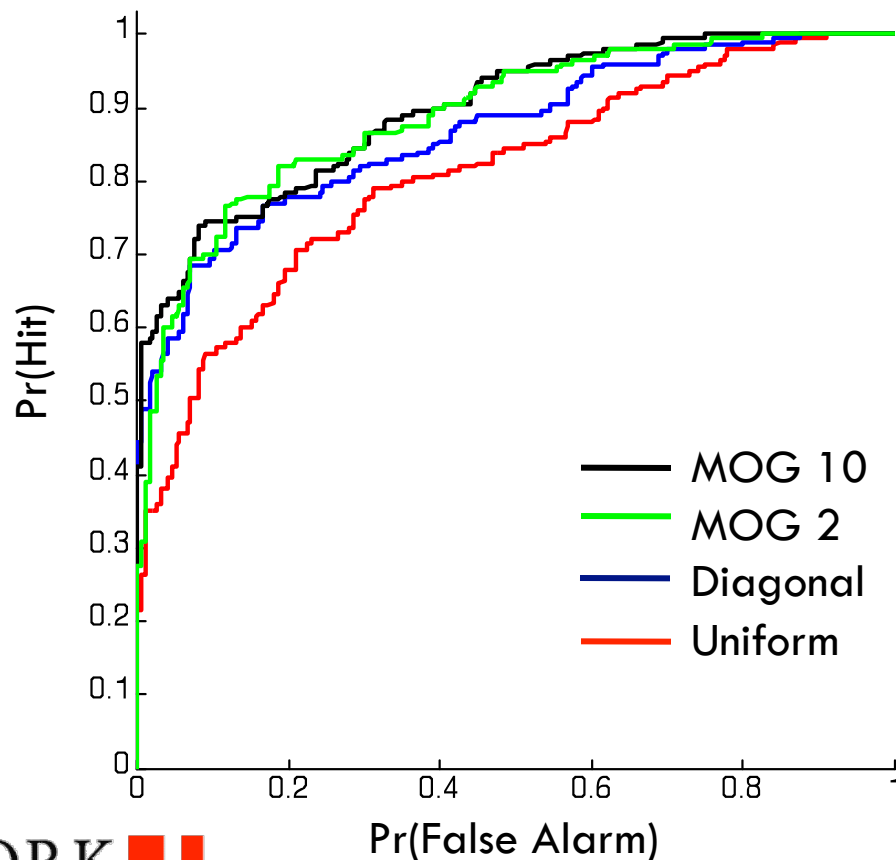
# MOG 10 Components



# Results for Mog 10 Model

27

Expectation Maximization



Performance improves slightly more, particularly at low false alarm rates.

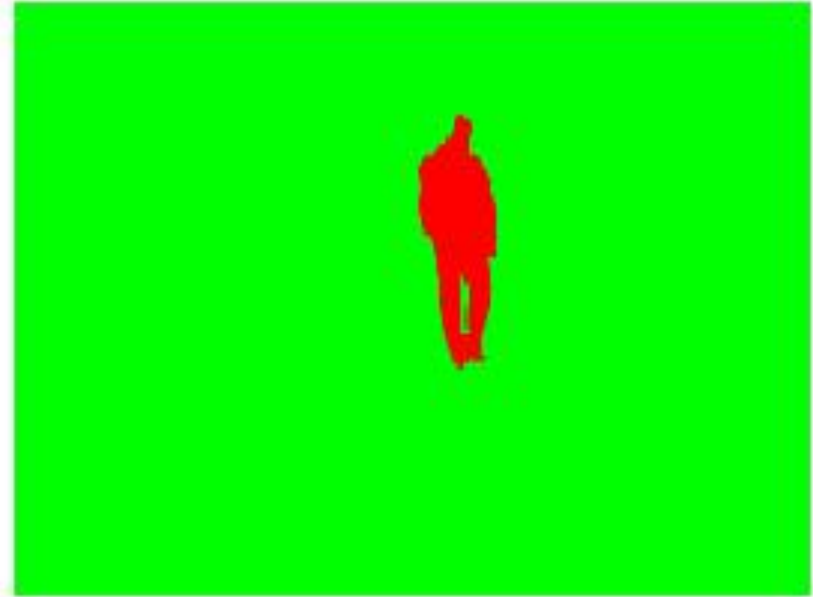
# Background Subtraction

28

Expectation Maximization



Test Image



Desired Segmentation

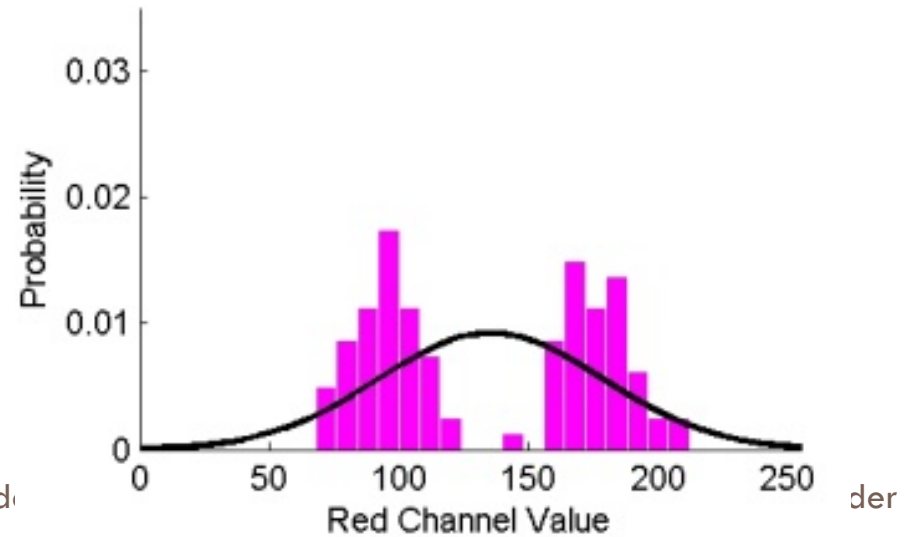
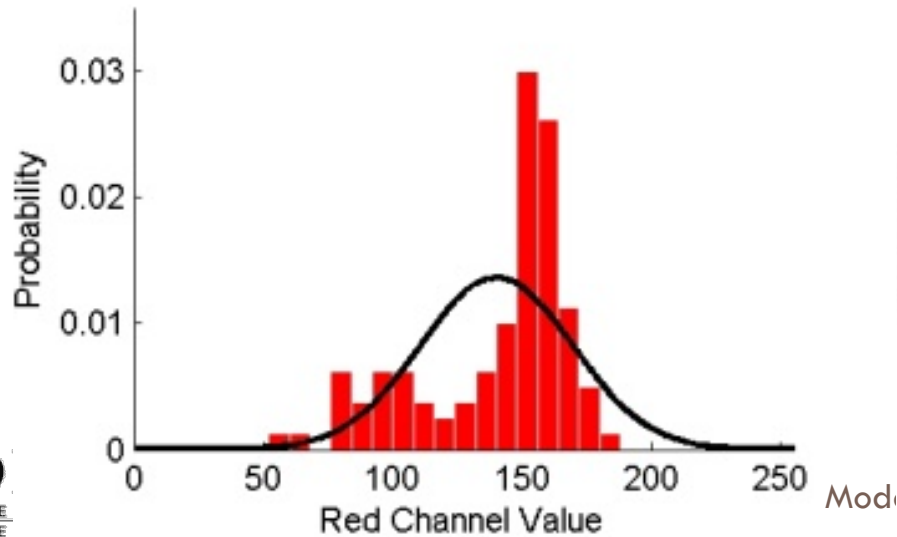
**GOAL :** (i) Learn background model (ii) use this to segment regions where the background has been occluded

# What if the scene isn't static?



Gaussian is no longer a good fit to the data.

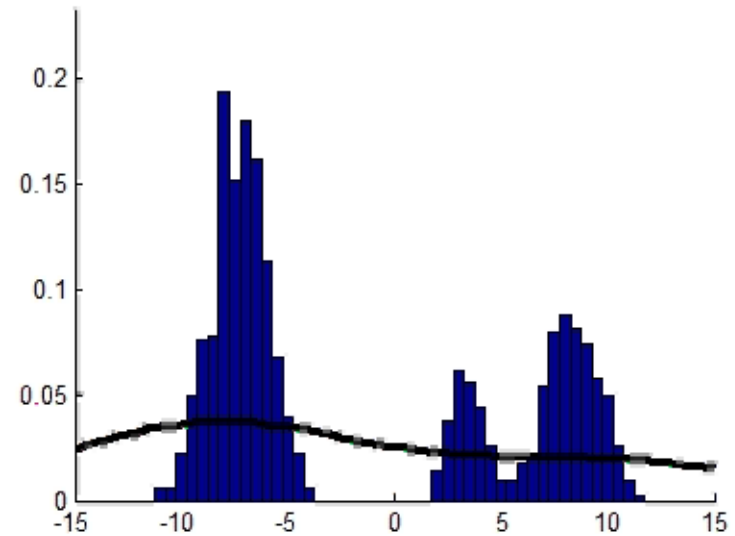
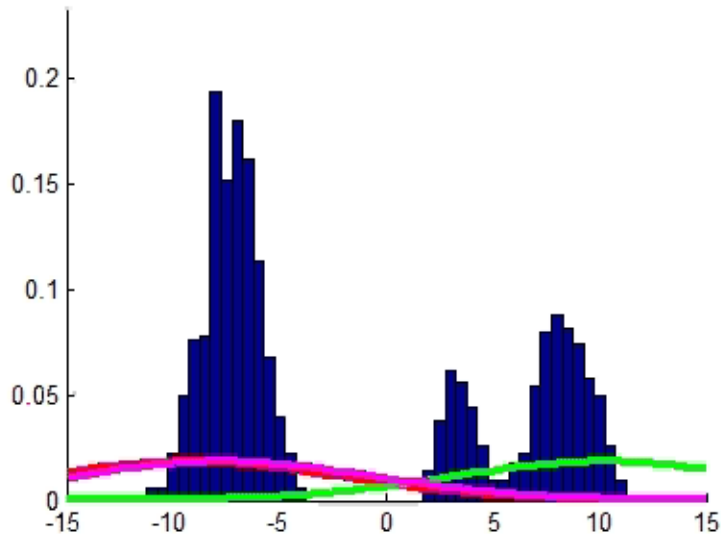
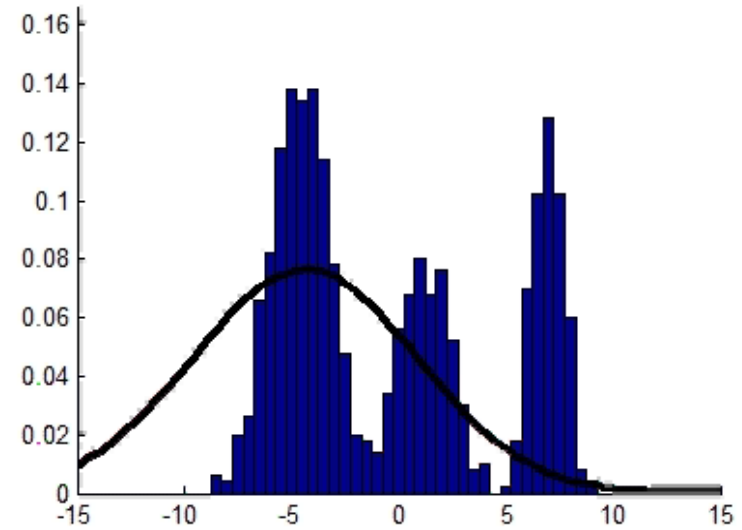
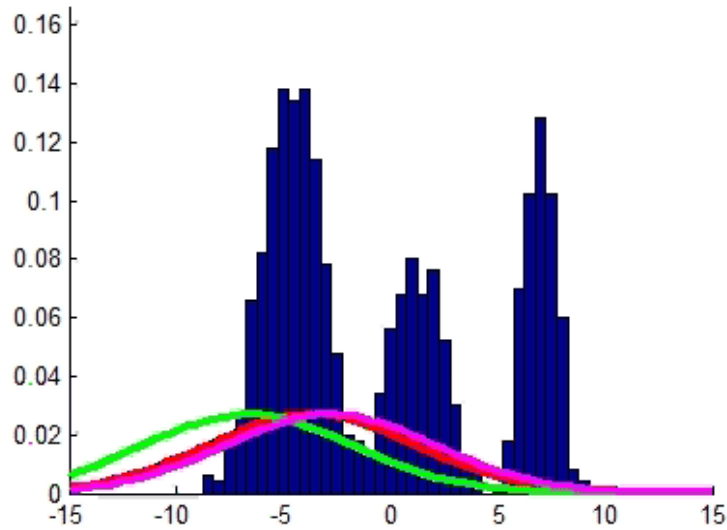
Not obvious exactly what probability model would fit better.



# One Dimensional Example

30

Expectation Maximization



# Final Fitted Models

31

Expectation Maximization

