

- Expectation Maximization:  
iterative algorithm for maximizing likelihood

- Given: vectors of "visible" variables  $v_n$   
Hidden: vectors of "hidden" variables  $h_n$   
Where  $n$  is example index

- Complete data set:  $\underbrace{U = \{v_n\}, H = \{h_n\}}_{\text{true data}}$

- Model specifies a joint distribution,

$$P(U, H | \theta)$$

$\uparrow$  parameters of model

- Usually i.i.d. data

$$P(U, H | \theta) = \prod P(v_n, h_n | \theta)$$

Then whole analysis "decomposes"!

$$- P(V | \theta) = \sum_H P(U, H | \theta)$$

$$= \sum_H P(V | H, \theta) P(H | \theta)$$

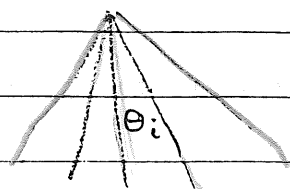
Goal:

- Maximize  $\ln P(V | \theta)$

Note: Its log of a sum over the hidden variables

## Example 1: Mixture of $m$ fixed densities

2



$P(x|i)$   
input distributions

$$\theta = \{\theta_1, \theta_2, \dots, \theta_m\}$$

Mixture coefficients

Statisticians often use  $\theta$  for their model / parameter

$$P(x|\theta) = \sum_i \theta_i P(x|i)$$

VISIBLE HIDDEN

$$\text{Formally: } P(x|\theta) = \sum_i P(x, i|\theta)$$

VISIBLE

PARAMETER

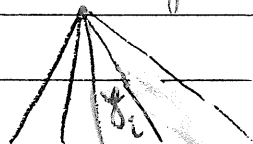
$$= \sum_i P(x|i, \theta) P(i|\theta)$$

fixed densities

$\theta_i$

## Example 2: Mixture of Gaussians

Mixture of Gaussians



$P(x|i, \theta)$

PARAMETERS

$$\theta = \{\gamma_i\} \cup \{\mu_i\} \cup \{\Sigma_i\} \quad 1 \leq i \leq k$$

mixture coefficients  
 $k$

means  
 $d \times 1$

co-variance matrices  
 $d \times d$

$$P(x|\theta) = \sum_i P(i|\theta) P(x|i, \theta)$$

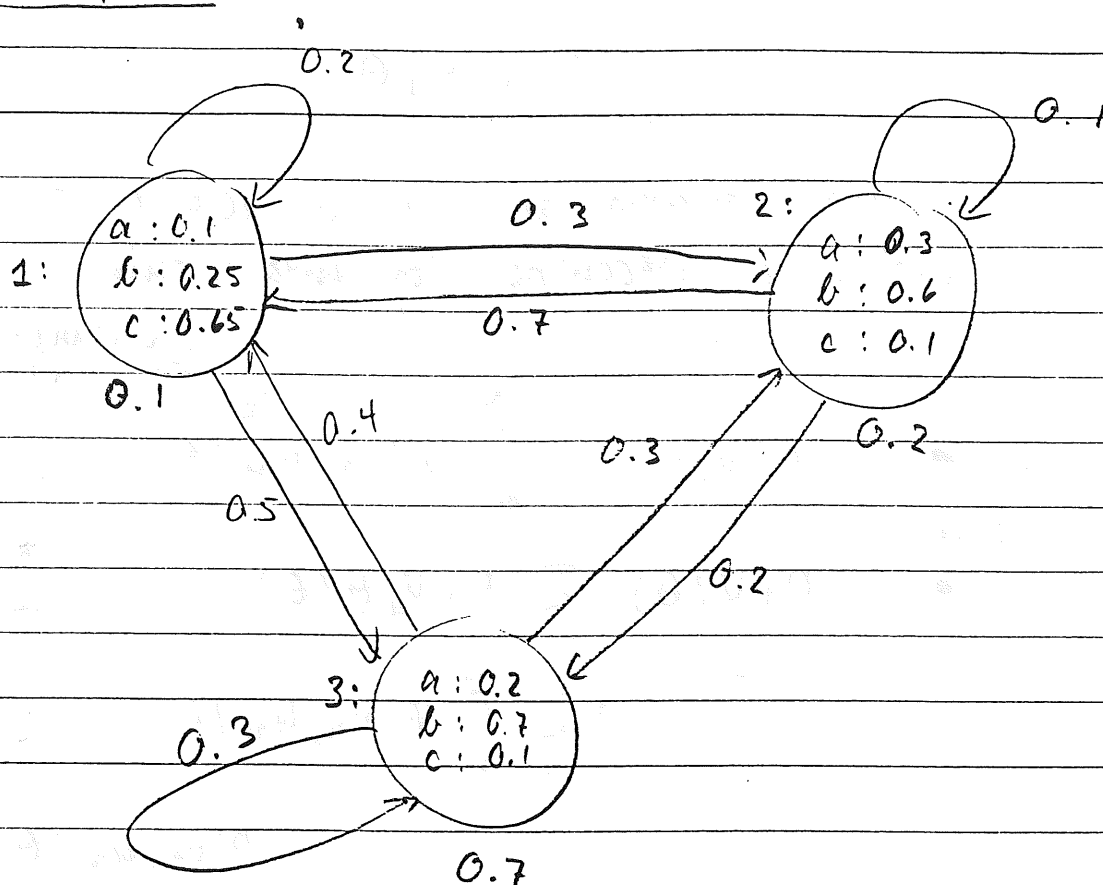
$$= \sum_i \gamma_i \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)}$$

POINT  $\bar{x} \in \mathbb{R}^d$  VISIBLE

$i \in \{1, 2, \dots, k\}$  HIDDEN

$\uparrow$  # OF GAUSSIANS

# Example 3: Hidden Markov Models



$$P(\underbrace{a b c a}_{\text{observed sequence}}, \underbrace{1 2 3 1}_{\text{hidden state sequence}} | \Theta) = 0.1 \cdot 0.1 \cdot 0.3 \cdot 0.6 \cdot 0.2 \cdot 0.1 \cdot 0.4 \cdot 0.1$$

1
2
3
1

parameters

PARAMETERS:  $\Theta = \{ \theta_i \}$  initial state output transition probabilities

$1 \leq i \leq m$

$$P(x, s | \Theta) = \prod_{i=1}^m \theta_i^{n_i(x, s)},$$

where  $n_i$  is # of times  $\theta_i$  occurs in  $(x, s)$

$$P(x | \Theta) = \sum_s P(x, s | \Theta)$$

sum over all hidden paths

Loss :  $-\ln P(V|\theta)$

$$= -\ln \sum_H P(V, H|\theta)$$

Hard to minimize because log of sum!

Iid. case: Decomposition into a sum

$$V = \{v_n\}, \quad H = \{h_n\}$$

SEQUENCES OF  
FIXED LENGTH  $N$  | ONE  
( $v_n, h_n$ )  
PER EXAMPLE

ASSUMPTION:  $P(V, H|\theta) = \prod_{n=1}^N P(v_n, h_n|\theta)$

THIS IMPLIES:

$$P(V|\theta) = \sum_H P(V, H|\theta)$$

$$= \sum_H \prod_n P(v_n, h_n|\theta)$$

$$= \sum_{h_1, h_2, \dots, h_N} \prod_n P(v_n, h_n|\theta)$$

$$= \prod_n \sum_{h_n} P(v_n, h_n|\theta)$$

$$= \prod_n P(v_n|\theta)$$

$$P(H|V, \theta) = \frac{P(H, V|\theta)}{P(V|\theta)} = \frac{\prod_n P(h_n, v_n|\theta)}{\prod_n P(v_n|\theta)} = \prod_n P(h_n|v_n, \theta)$$

$$-\ln P(V|\theta) = -\sum_n \ln P(v_n|\theta)$$

STILL SUM:  $\sum_{h_n} P(v_n, h_n|\theta)$

Product becomes sum

Loss is essentially a relative entropy

$$\sum_n \frac{1}{N} \ln \frac{1}{P(v_n|\theta)} = -\frac{1}{N} \sum_n \ln P(v_n|\theta) - \ln N$$

T.I.D. MAKES RELATIVE ENTROPIES  
DECOMPOSE AS WELL

$$\begin{aligned}
 & \sum_H P(H|U, \theta) \ln \frac{P(H|U, \theta)}{P(H|U, \tilde{\theta})} \\
 &= \sum_{\{h_n\}} \prod_n P(h_n|v_n, \theta) \ln \frac{\prod_n P(h_n|v_n, \theta)}{\prod_n P(h_n|v_n, \tilde{\theta})} \\
 &= \sum_n \sum_{h_n} P(h_n|v_n, \theta) \ln \frac{P(h_n|v_n, \theta)}{P(h_n|v_n, \tilde{\theta})}
 \end{aligned}$$

SIMPLEST CASE

$$\begin{aligned}
 & \sum_{x,y} P(x) P(y) \ln \frac{P(x) P(y)}{q(x) q(y)} \\
 &= \sum_{x,y} P(x) P(y) \left( \ln \frac{P(x)}{q(x)} + \ln \frac{P(y)}{q(y)} \right) \\
 &= \sum_x P(x) \ln \frac{P(x)}{q(x)} + \sum_y P(y) \ln \frac{P(y)}{q(y)}
 \end{aligned}$$

UPSHOT:

ALL BECOMES SUMS OVER EXAMPLES

$$\sum_H P(H|V, \theta) \ln \frac{P(H|V, \theta)}{P(H|V, \tilde{\theta})} + \eta \underbrace{(-\ln P(V|\tilde{\theta}))}_b \quad (*)$$

Divergence that motivates EM      Too hard to minimize!

$$= \sum_H P(H|V, \theta) \ln \frac{P(H, V|\theta)}{P(H, V|\tilde{\theta})} - \eta \ln P(V|\tilde{\theta})$$

$$= \sum_H P(H|V, \theta) \ln \frac{P(H, V|\theta)}{P(H, V|\tilde{\theta})} + \cancel{\ln P(V|\tilde{\theta})} - \ln P(V|\theta) - \eta \ln P(V|\tilde{\theta})$$

$$\eta = 1$$

$$= \sum_H P(H|V, \theta) \ln \frac{P(H, V|\theta)}{P(H, V|\tilde{\theta})} - \ln P(V|\theta)$$

easy to minimize

constant

Estimation Step:

- Compute posterior  $P(H|V, \theta)$

Maximization Step:

$$\tilde{\theta} := \operatorname{argmax}_{\tilde{\theta}} \sum_H P(H|V, \theta) \ln P(H, V|\tilde{\theta})$$

IID  
CASE:

$$P(V, H|\theta) = \prod_{n=1}^N P(v_n, h_n|\theta)$$

$$(*) \text{ becomes: } - \sum_n \sum_{h_n} P(h_n|v_n, \theta) \ln P(h_n, v_n|\tilde{\theta}) + \text{constant}$$

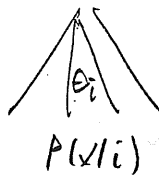
E-step:

Compute posteriors  $P(h_n|v_n, \theta)$

M-step:

$$\tilde{\theta} = \operatorname{argmax}_{\tilde{\theta}} \sum_n \sum_{h_n} P(h_n|v_n, \theta) \ln P(h_n, v_n|\tilde{\theta})$$

Example 1:



Parameters  $\theta = \{\theta_1, \dots, \theta_m\}$

visible  $x$

hidden  $i$

fixed densities:  $P(x|i, \theta) = P(x|i)$

$$P(x|\theta) = \sum_i P(x, i|\theta)$$

$$= \sum_i \underbrace{P(i|\theta)} \underbrace{P(x|i, \theta)}$$

$$= \sum_i \theta_i P(x|i)$$

E-step:

$$P(i|x_n, \theta) = \frac{P(i, x_n|\theta)}{P(x_n|\theta)} = \frac{\theta_i P(x_n|i)}{\sum_j \theta_j P(x_n|j)}$$

M-step:

EXPECTED USAGE  
OF DISTRIBUTION  $i$

Maximize

$$\sum_n \sum_{x_n} P(i|x_n, \theta) \ln P(x_n, i|\tilde{\theta}) + \lambda (\sum \tilde{\theta}_i - 1)$$

$$= \sum_n \sum_{x_n} P(i|x_n, \theta) \ln \underbrace{\tilde{\theta}_i P(x_n|i)}_{\text{constant}} + \lambda (\sum \tilde{\theta}_i - 1)$$

$$\frac{\partial}{\partial \tilde{\theta}_i} = \sum_n \frac{P(i|x_n, \theta)}{\tilde{\theta}_i} + \lambda = 0$$

$$\Leftrightarrow \sum_n P(i|x_n, \theta) + \lambda \tilde{\theta}_i = 0 \quad (*)$$

$$\sum_i \sum_n P(i|x_n, \theta) + \lambda \sum_i \tilde{\theta}_i = 0$$

$$\lambda = -N$$

$$\text{From } (*): \tilde{\theta}_i = \frac{1}{N} \sum_n P(i|x_n, \theta)$$

Average Posterior

$N=1$  Bayes Rule

## Example 2: Mixture of Gaussians

$P(x|i, \theta)$  not independent of  $\theta$

$$\begin{aligned} P(x|\theta) &= \sum_i \underbrace{P(i|\theta)}_{\gamma_i} \underbrace{P(x|i, \theta)}_{\frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)}} \\ &= \sum_i \gamma_i \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1} (x-\mu_i)} \end{aligned}$$

E step:

$$P(i|x_n, \theta) = \frac{P(i, x_n|\theta)}{P(x_n|\theta)} = \frac{\gamma_i P(x|i, \theta)}{\sum_j \gamma_j P(x|j, \theta)}$$

M step:

Maximize

$$\begin{aligned} &\sum_n \sum_i P(i|x_n, \theta) \ln P(i, x_n|\tilde{\theta}) \\ &= \sum_n \sum_i P(i|x_n, \theta) \ln \underbrace{P(i|\tilde{\theta})}_{\gamma_i} P(x_n|i, \tilde{\theta}) \quad (*) \end{aligned}$$

- 3 classes of parameters:
- a) mixture coefficients  $\gamma_i$
  - b) means  $\mu_i$
  - c) covariance matrices  $\Sigma_i$

Maximize above for one class at a time while keeping the other classes fixed.



a)  $\mu_i$  and  $\Sigma_i$  fixed

Thus  $P(x_n | i, \tilde{\theta}) = P(x_n | i, \theta)$

(\*) becomes  $\sum_n \sum_i P(i | x_n, \theta) \ln \tilde{\gamma}_i + \text{constant}$

Minimize as before with Lagrangian:

$$\tilde{\gamma}_i = \frac{1}{N} \sum_n P(i | x_n, \theta)$$

b) and c):  $\gamma_i$  fixed

Thus  $P(i | \tilde{\theta}) = P(i | \theta) = \gamma_i$

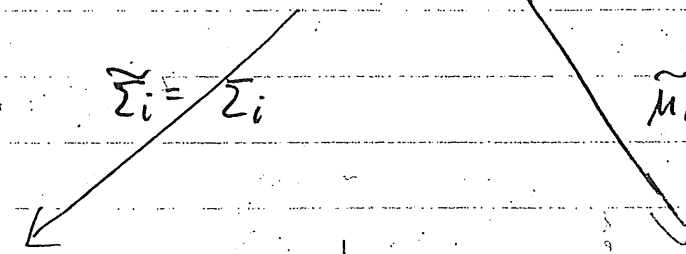
(\*) becomes  $\sum_n \sum_i P(i | x_n, \theta) \ln P(x_n | i, \tilde{\theta}) + \text{const.}$

$$\sum_n \sum_i P(i | x_n, \theta) \ln \left( (2\pi)^{-d/2} |\tilde{\Sigma}_i|^{-1/2} e^{-\frac{1}{2} (x_n - \tilde{\mu}_i)^T \tilde{\Sigma}_i^{-1} (x_n - \tilde{\mu}_i)} \right) + \text{const.}$$

$$\sum_n \sum_i P(i | x_n, \theta) \left( -\frac{1}{2} \ln |\tilde{\Sigma}_i| - \frac{1}{2} (x_n - \tilde{\mu}_i)^T \tilde{\Sigma}_i^{-1} (x_n - \tilde{\mu}_i) \right) + \text{const}$$

$$\tilde{\Sigma}_i = \Sigma_i$$

$$\tilde{\mu}_i = \mu_i$$



$$\sum_n \sum_i P(i | x_n, \theta) \left( -\frac{1}{2} (x_n - \tilde{\mu}_i)^T \tilde{\Sigma}_i^{-1} (x_n - \tilde{\mu}_i) \right) + \text{const} \quad \left| \quad \sum_n \sum_i P(i | x_n, \theta) \left( -\frac{1}{2} \ln |\tilde{\Sigma}_i| - \frac{1}{2} (x_n - \tilde{\mu}_i)^T \tilde{\Sigma}_i^{-1} (x_n - \tilde{\mu}_i) \right) + \text{const} \right.$$

b)

c)

$$b) \frac{\partial}{\partial \tilde{\mu}_i} = \sum_n P(i|x_n, \theta) \tilde{\Sigma}_i^{-1} (x_n - \tilde{\mu}_i) = 0$$

$$\Leftrightarrow \sum_n P(i|x_n, \theta) (x_n - \tilde{\mu}_i) = 0$$

$$\sum_n P(i|x_n, \theta) x_n = \sum_n P(i|x_n, \theta) \tilde{\mu}_i$$

$$\tilde{\mu}_i = \frac{\sum_n P(i|x_n, \theta) x_n}{\sum_n P(i|x_n, \theta)}$$

$$c) \frac{\partial}{\partial \tilde{\Sigma}_i} = \sum_n P(i|x_n, \theta) \left( -\frac{1}{2} \tilde{\Sigma}_i^{-1} + \frac{1}{2} \tilde{\Sigma}_i^{-1} (x_n - \mu_i) (x_n - \mu_i)^T \tilde{\Sigma}_i^{-1} \right) = 0$$

$$\Leftrightarrow \sum_n P(i|x_n, \theta) \left( -\frac{1}{2} \tilde{\Sigma}_i + \frac{1}{2} (x_n - \mu_i) (x_n - \mu_i)^T \right) = 0$$

$$\tilde{\Sigma}_i = \frac{\sum_n P(i|x_n, \theta) (x_n - \mu_i) (x_n - \mu_i)^T}{\sum_n P(i|x_n, \theta)}$$

High-level intuition:

In EM parameters are often updated to their usage when generating the data.

Intuition why EM works:

A) Minimizing  $-\ln P(V|\tilde{\theta}) = -\ln \sum_H P(V, H|\tilde{\theta})$

is hard because "ln" of a sum.

Minimizing  $-\sum_H P(H|V, \theta) \ln P(H, V|\tilde{\theta})$

"easy" when  $P(H, V|\tilde{\theta})$  is product!

EM avoids minimizing "ln" of a sum directly.

By adding distance, minimization simplified,  
 $-\ln P(V|\tilde{\theta})$

B)  $= -\ln \sum_H P(V, H|\tilde{\theta})$  has many symmetries.

For example by renaming any local minima

for a mixture of  $m$  distinct Gaussians becomes

$m!$  local minima.

In  $-\sum_H P(H|V, \theta) \ln P(H, V|\tilde{\theta})$  not as many

symmetries. Hidden variables are "coupled"

with visible variables.

Example 3: HMM's  $P(x, s | \theta) = \prod_{i=1}^I \theta_i^{n_i(x, s)}$

$$\begin{aligned} \text{Maximize } & \sum_n \sum_{s_n} P(s_n | x_n, \theta) \ln P(x_n, s_n | \tilde{\theta}) \\ &= \sum_n \sum_{s_n} P(s_n | x_n, \theta) \sum_i n_i(x_n, s_n) \ln \tilde{\theta}_i \\ &= \sum_n \sum_i \ln \tilde{\theta}_i \underbrace{\sum_{s_n} P(s_n | x_n, \theta) n_i(x_n, s_n)}_{\hat{n}_i(x | \theta)} \end{aligned}$$

$\hat{n}_i(x | \theta)$   
Expected "usage" of param.  $\theta_i$   
Computable by dynamic program.

$[i] := \{j : \theta_i \text{ and } \theta_j \text{ in same "class"}\}$

All parameters of a class must sum to one.

- two classes associated with a state.
- one class associated with the initial state probabilities

Maximize

$$\sum_n \sum_i (\ln \tilde{\theta}_i) \hat{n}_i(x_n | \theta) + \sum_{[i]} \lambda_{[i]} \left( \sum_{j \in [i]} \theta_j - 1 \right)$$

not right

$$\frac{\partial}{\partial \tilde{\theta}_i} = \sum_n \frac{\hat{n}_i(x_n | \theta)}{\tilde{\theta}_i} + \lambda_{[i]} = 0$$

$$\tilde{\theta}_i = \frac{\sum_n \hat{n}_i(x_n | \theta)}{\lambda_{[i]}}$$

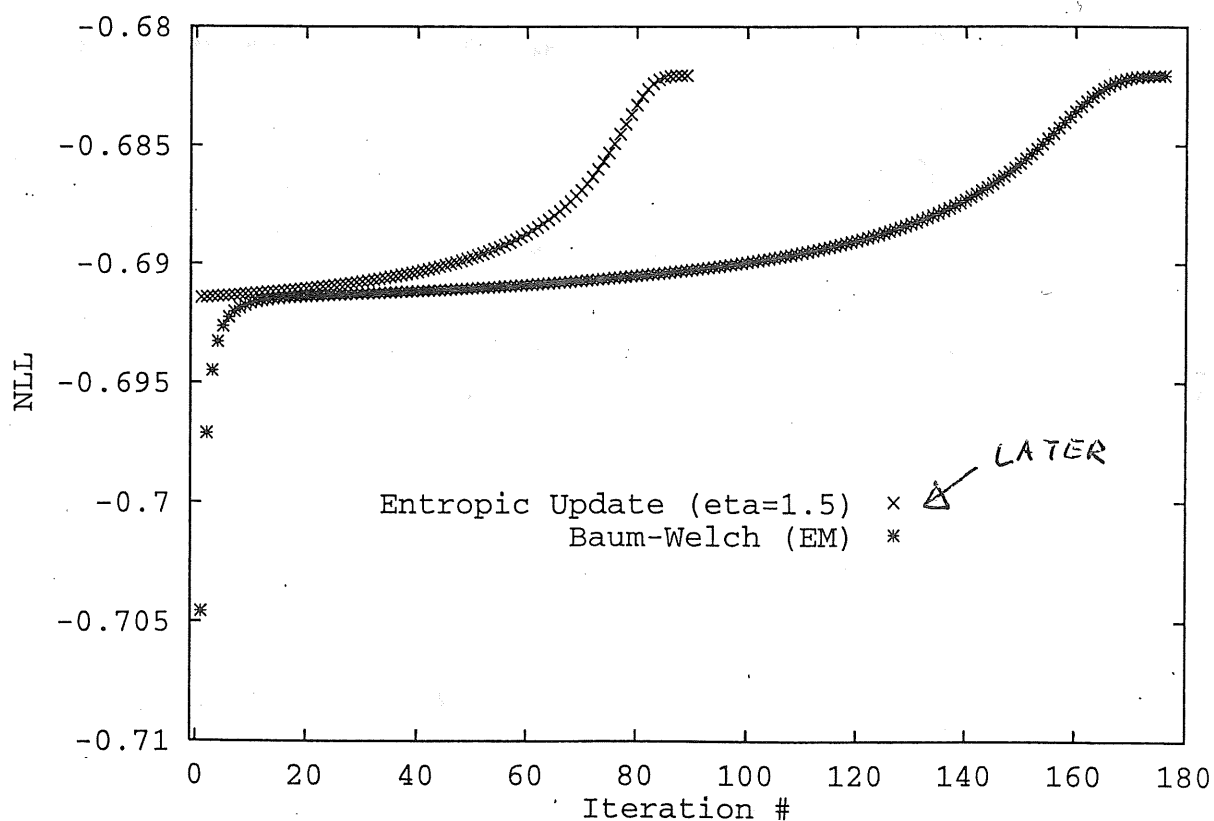
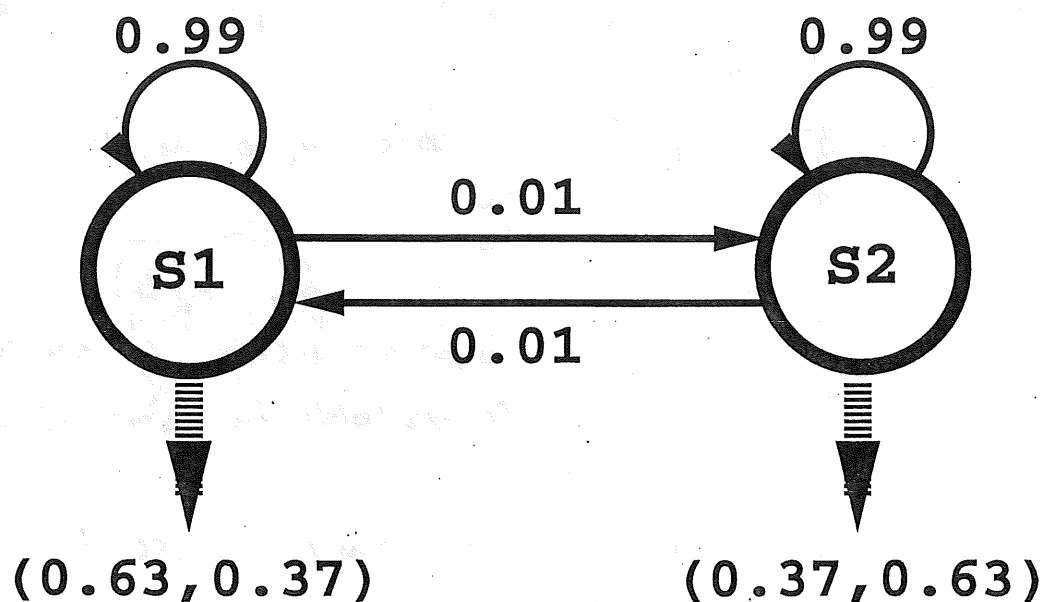
Enforcing constraint:  $\sum_{j \in [i]} \tilde{\theta}_j = 1$

$$\tilde{\theta}_i = \frac{\sum_n \hat{n}_i(x_n | \theta)}{\sum_{j \in [i]} \sum_n \hat{n}_j(x_n | \theta)}$$

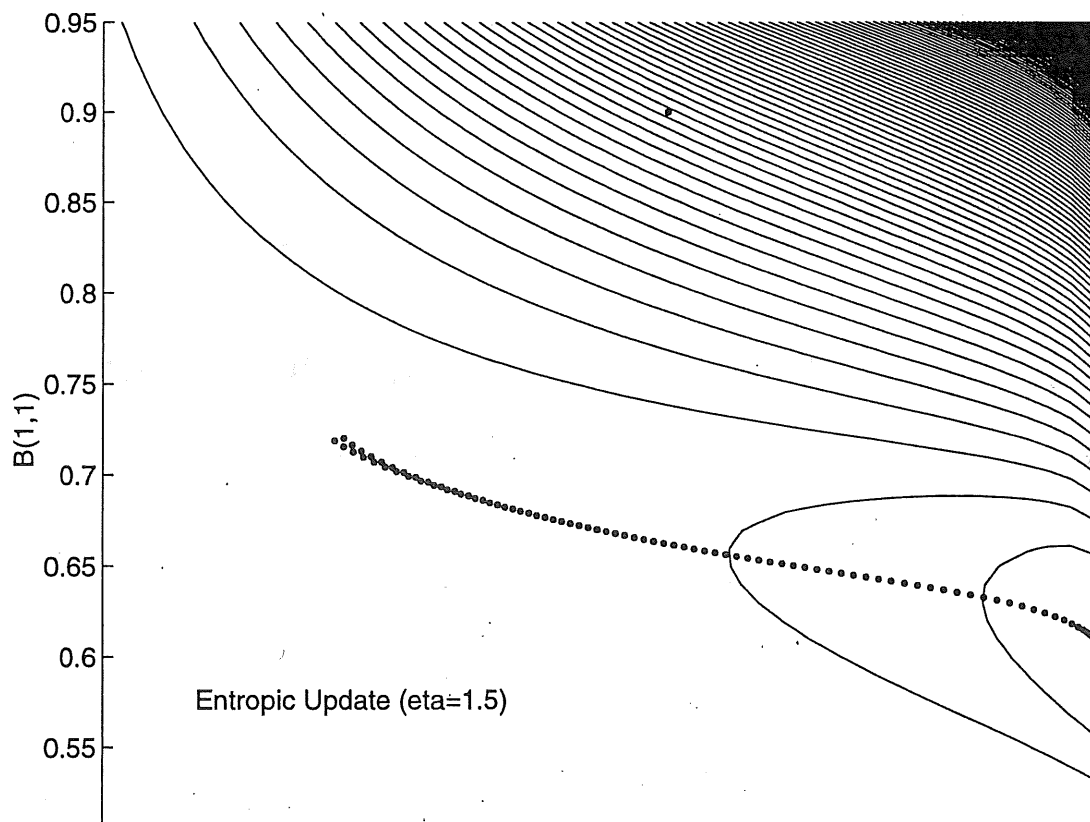
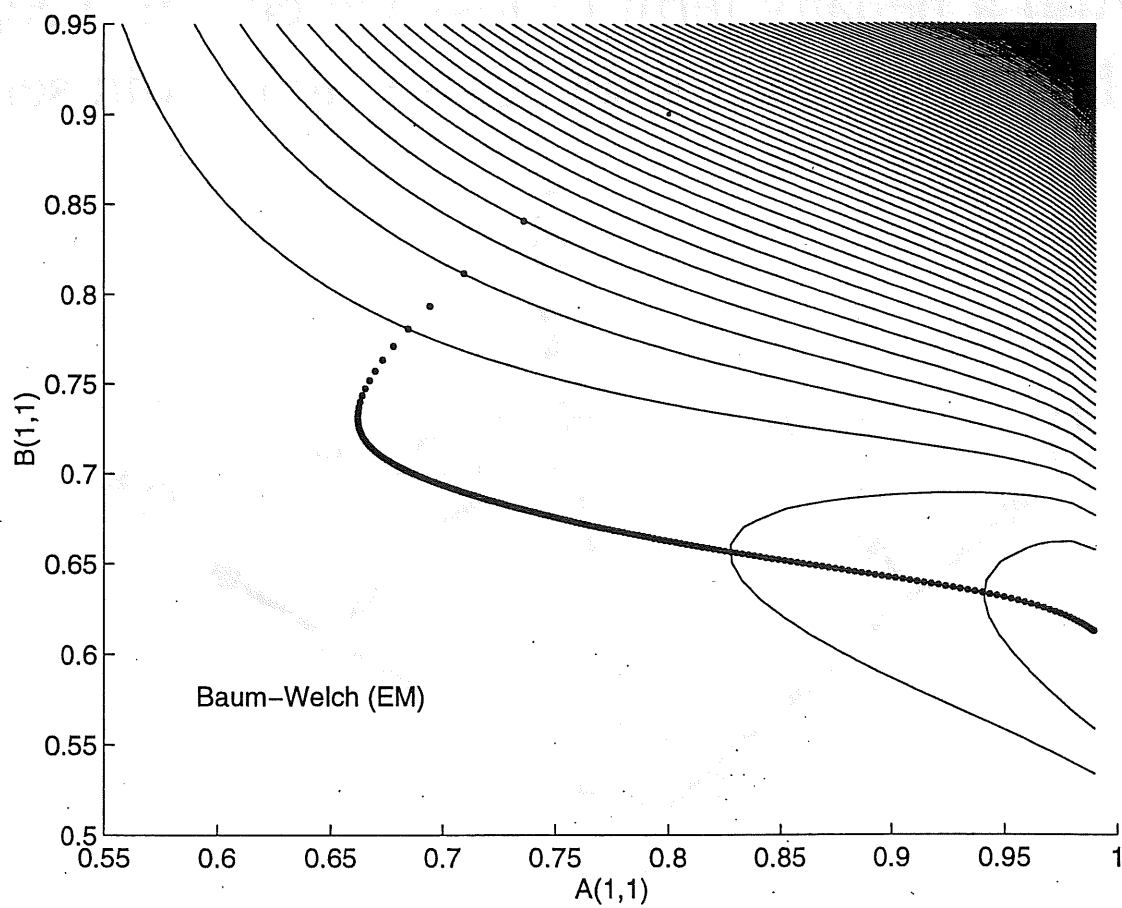
EM is too slow!

13

## Synthetic Data 1: 2 State HMM

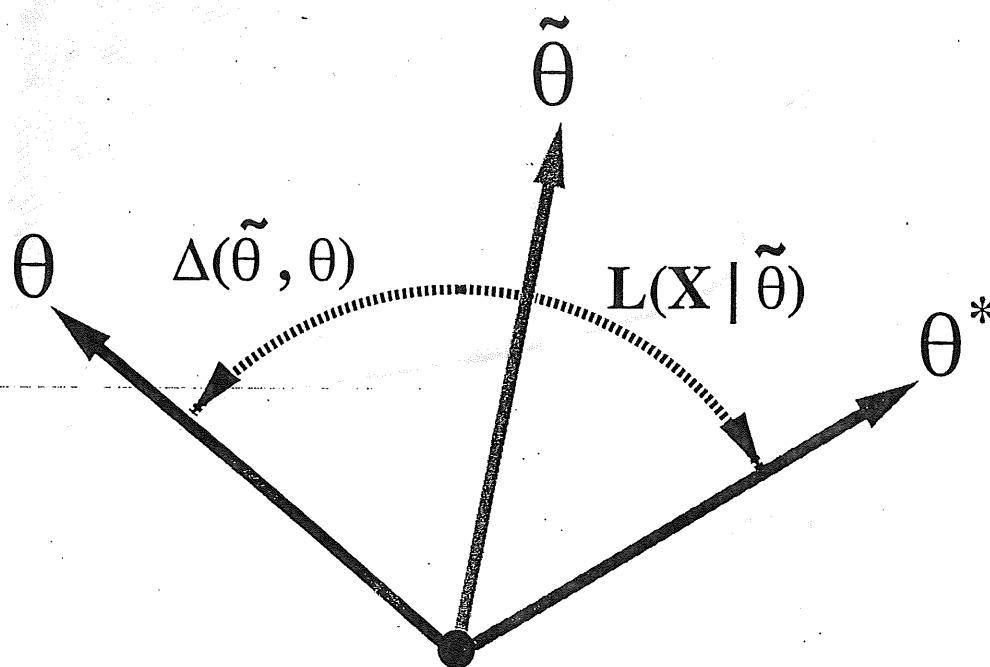


## Synthetic Data 1 (cont.)



## Framework for Parameter Update

- Add a penalty term to loss which will keep the new parameter set “close” to the old set.



- Solve:

$$\begin{aligned}\Theta^{t+1} &= \arg \min_{\tilde{\Theta}} U^t(\tilde{\Theta}) \\ U^t(\tilde{\Theta}) &= \Delta(\tilde{\Theta}, \Theta^t) + \eta \text{loss}(S | \tilde{\Theta})\end{aligned}\quad (*)$$

$\eta$  is a non-negative trade-off parameter that becomes the learning rate of the algorithm.

Any update of the form (\*) called

Implicit Update

# Minimal Properties of the Divergence

- (i)  $\Delta(\Theta, \Theta) = 0$
- (ii)  $\Delta(\tilde{\Theta}, \Theta) > 0$  whenever  $\tilde{\Theta} \neq \Theta$ .

## Key Lemma

If  $U^t(\tilde{\Theta}) < U^t(\Theta^t)$   
then  $\text{loss}(S|\tilde{\Theta}) < \text{loss}(S|\Theta^t)$ .

**Proof:**

$$\begin{aligned} U^t(\tilde{\Theta}) &= \Delta(\tilde{\Theta}, \Theta^t) + \eta \text{loss}(S|\tilde{\Theta}) \\ &< U^t(\Theta^t) = \Delta(\Theta^t, \Theta^t) + \eta \text{loss}(S|\Theta^t) \stackrel{(i)}{=} \eta \text{loss}(S|\Theta^t) \end{aligned}$$

This is equivalent to

$$\begin{aligned} \text{loss}(S|\tilde{\Theta}) &= \text{loss}(S|\Theta^t) - \frac{\Delta(\tilde{\Theta}, \Theta^t)}{\eta} \\ &\stackrel{(ii)}{<} \text{loss}(S|\Theta^t) \end{aligned}$$

For any implicit update s.t.  $\Theta^{t+1} \neq \Theta^t$ ,

$$\text{loss}(S|\Theta^{t+1}) < \text{loss}(S|\Theta^t)$$



What's good about EM:

- implicit update. Thus negative log likelihood decreases in each iteration
- simple and elegant

Bad: Converges too slowly

$$\sum_H P(H|V, \theta) \ln \frac{P(H|V, \theta)}{P(H|V, \tilde{\theta})} - \eta \ln P(V|\tilde{\theta})$$

datadata

Simplification for  $\eta = 1$



$$- \sum_H P(H|V, \theta) \ln P(H, V|\tilde{\theta}) + \text{const.}$$

Want  $\eta > 1$ . In that case simplification does not work!!!

Idea 1: Use  $\eta > 1$  but approximate  $-\ln P(V|\tilde{\theta})$

by 1. order Taylor.

Does not seem to work.

Idea 2: Use  $\eta > 1$ , different distance, and 1. order Taylor.

$$\sum_{V'} \sum_H P(H, V'|\tilde{\theta}) \ln \frac{P(H, V'|\tilde{\theta})}{P(H, V'|\theta)} - \eta (\ln P(V|\tilde{\theta}) + (\tilde{\theta} - \theta) \frac{\partial \ln P(V|\theta)}{\partial \theta})$$

data

Explicit!

Distance we use:

18

- different direction of entropy
- $V$  integrated over domain
- avoids "ln" of sum in different way.

In all three examples our method converges faster.

$$\left. \frac{\partial -\ln P(V|\tilde{\Theta}_\eta)}{\partial \eta} \right|_{\eta=0} < 0 \quad \text{unless at extremum}$$

Provided  $\eta$  is close enough to 0,  
then loss decreases.

We don't know why our method is so good.