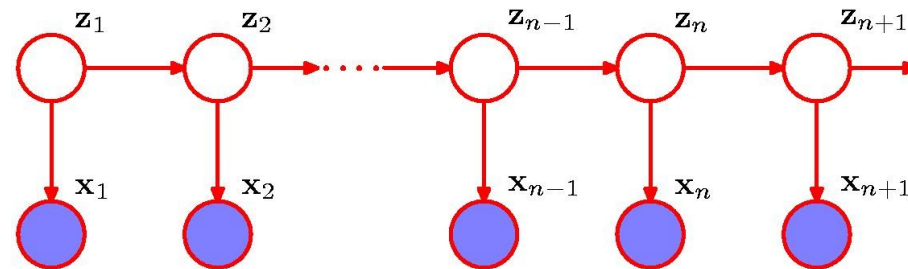


# Hidden Markov Models

## Training – Selecting model parameters



Christian Nørgaard Storm Pedersen

[cstorm@birc.au.dk](mailto:cstorm@birc.au.dk)

## What we know

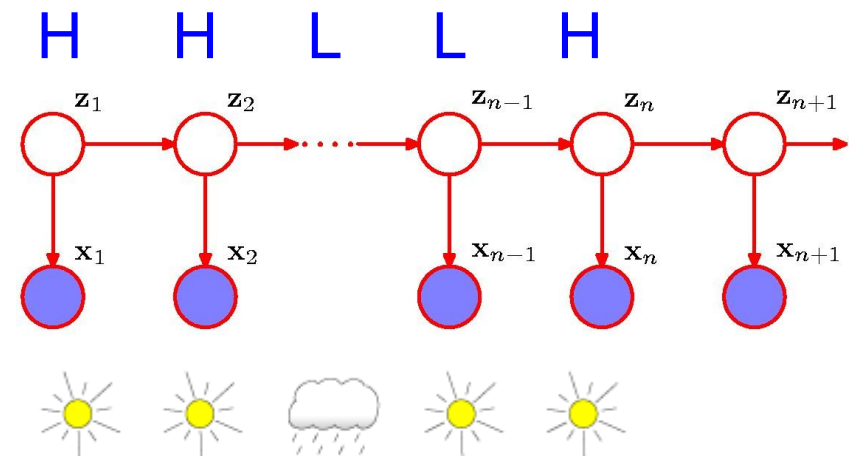
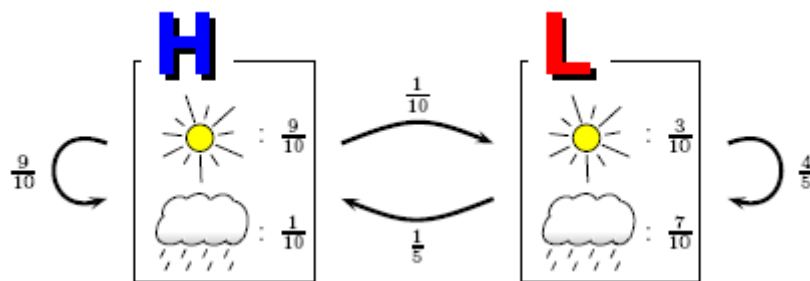
- The terminology and notation of hidden Markov models (**HMMs**)
- The **forward- and backward-algorithms** for determining the likelihood  $p(\mathbf{X})$  of a sequence of observations, and computing the **posterior decoding**.
- The **Viterbi-algorithm** for finding the most likely explanation (sequence of latent states) of a sequence of observations.
- How to implement the Viterbi-algorithm using log-transform (and the forward- and backward-algorithms using scaling).

## Now

- Training, or how to select model parameters (transition and emission probabilities) to reflect either a set of corresponding  $(\mathbf{X}, \mathbf{Z})$ 's, (or just a set of  $\mathbf{X}$ 's) ...

# Selecting “the right” parameters

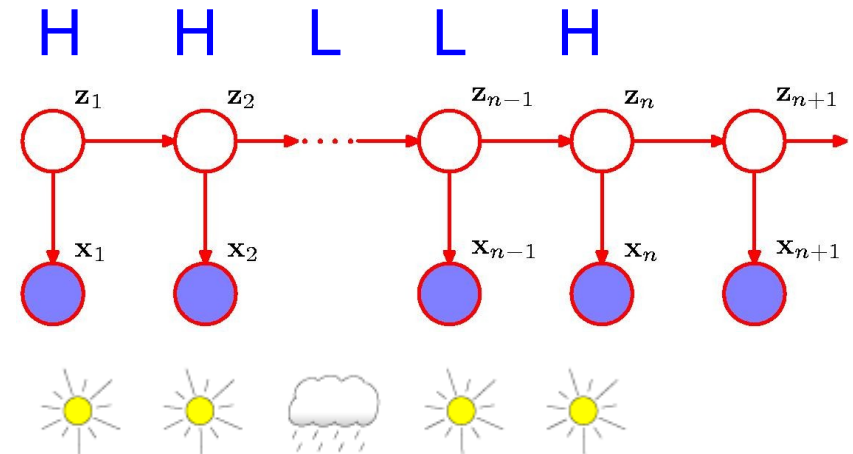
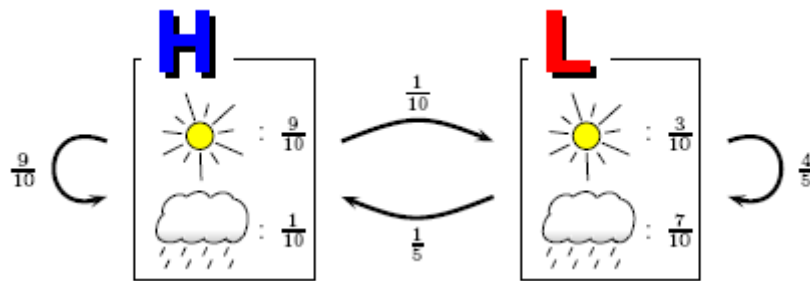
Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...



How should we set the model parameters, i.e. transition  $\mathbf{A}$ ,  $\boldsymbol{\pi}$ , and emission probabilities  $\boldsymbol{\Phi}$ , to make the given  $(\mathbf{X}, \mathbf{Z})$ 's most likely?

# Selecting “the right” parameters

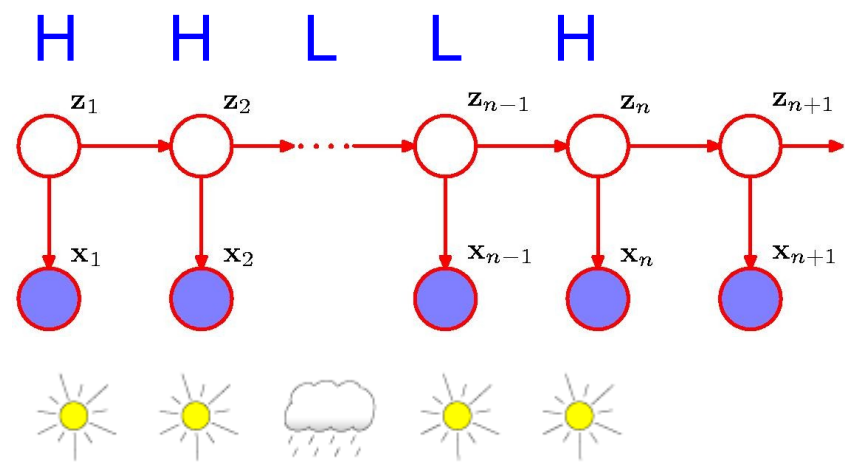
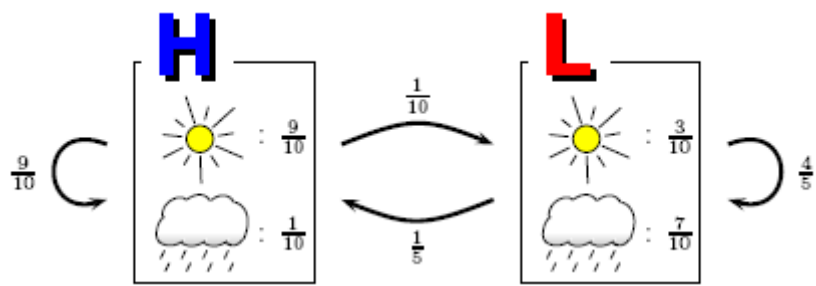
Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...



How should we set the model parameters, i.e. transition  $\mathbf{A}$ ,  $\boldsymbol{\pi}$ , and emission probabilities  $\boldsymbol{\Phi}$ , to make the given  $(\mathbf{X}, \mathbf{Z})$ 's most likely?

**Intuition:** The parameters should reflect what we have seen ...

# Selecting “the right” transition probs



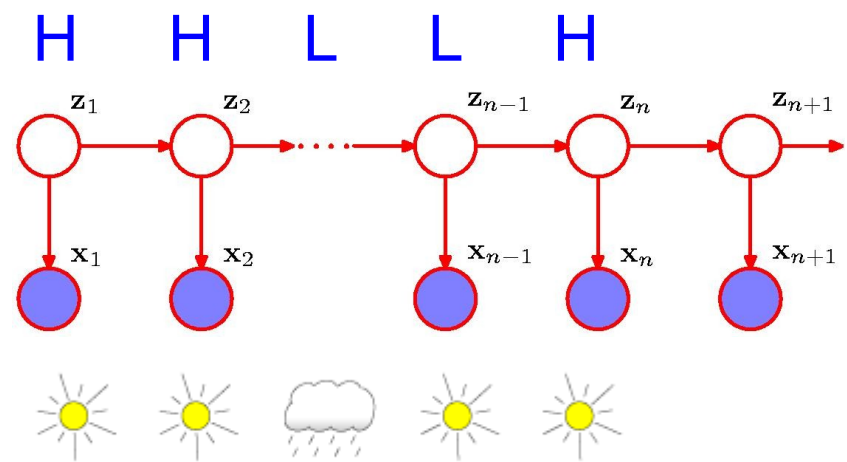
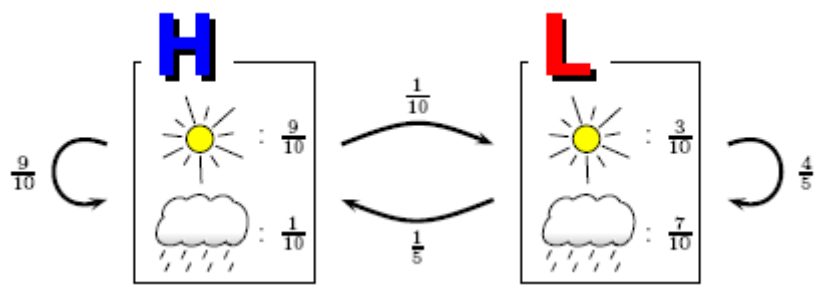
$A_{jk}$  is the probability of a transition from state  $j$  to state  $k$ , and  $\pi_k$  is the probability of starting in state  $k$  ...

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} =$$

How many times is the transition from state  $j$  to state  $k$  taken

How many times is a transition from state  $j$  to any state taken

# Selecting “the right” transition probs



$A_{jk}$  is the probability of a transition from state  $j$  to state  $k$ , and  $\pi_k$  is the probability of starting in state  $k$  ...

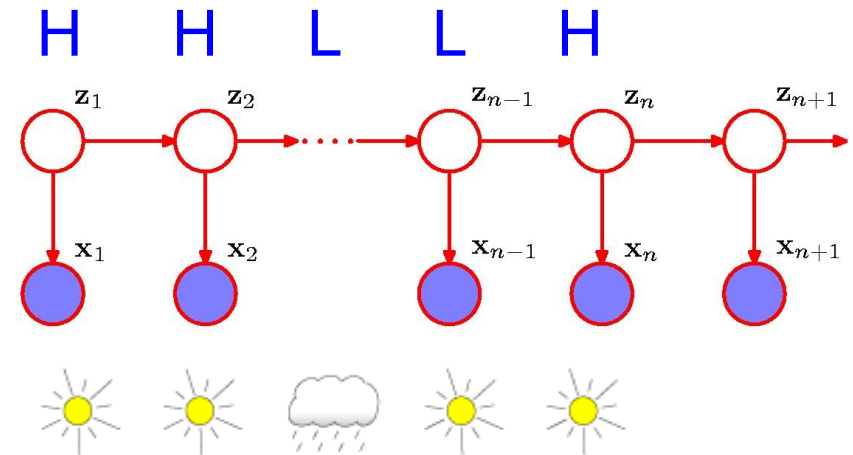
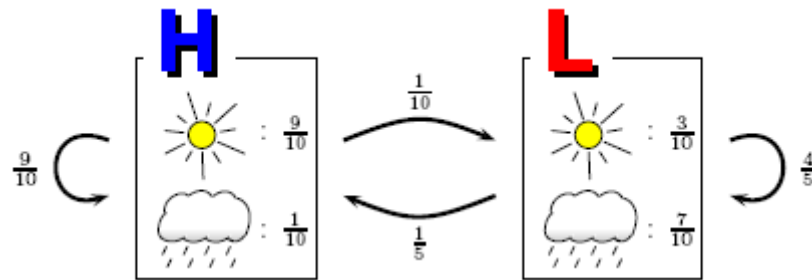
$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} =$$

$$\pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}}$$

How many times is the transition from state  $j$  to state  $k$  taken

How many times is a transition from state  $j$  to any state taken

# Selecting “the right” emission probs



If we assume discrete observations, then  $\phi_{ik}$  is the probability of emitting symbol  $i$  from state  $k$  ...

$$\phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}} =$$

How many times is symbol  $i$  emitted from state  $k$

How many times is a symbol emitted from state  $k$

# Selecting “the right” parameters

Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

We simply count how many times each outcome of the multinomial variables (a transition or emission) is observed ...



# Selecting “the right” parameters

Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

We simply count how many times each outcome of the multinomial variables (a transition or emission) is observed ...

This yield a maximum likelihood estimate (MLE)  $\boldsymbol{\theta}^*$  of  $p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta})$ , which is what we mathematically want ...

$$f_{\mathbf{X}, \mathbf{Z}}(\boldsymbol{\theta}) = p(\mathbf{X}, \mathbf{Z} \mid \boldsymbol{\theta}) \quad \boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} f_{\mathbf{X}, \mathbf{Z}}(\boldsymbol{\theta})$$

# Selecting “the right” parameters

Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

We simply count how many times each outcome of the multinomial variables (a transition or emission) is observed ...

This yield a maximum likelihood estimate (MLE)  $\boldsymbol{\theta}^*$  of  $p(\mathbf{X},\mathbf{Z} \mid \boldsymbol{\theta})$ , which is what we mathematically want ...

**Any problems?**

# Selecting “the right” parameters

Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

We simply count how many times each outcome of the multinomial variables (a transition or emission) is observed ...

This yield a maximum likelihood estimate (MLE)  $\boldsymbol{\theta}^*$  of  $p(\mathbf{X},\mathbf{Z} \mid \boldsymbol{\theta})$ , which is what we mathematically want ...

**Any problems?** What if e.g. the transition from state  $j$  to  $k$  is *not* observed, then probability  $A_{jk}$  is set to 0.

# Selecting “the right” parameters

Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are given ...

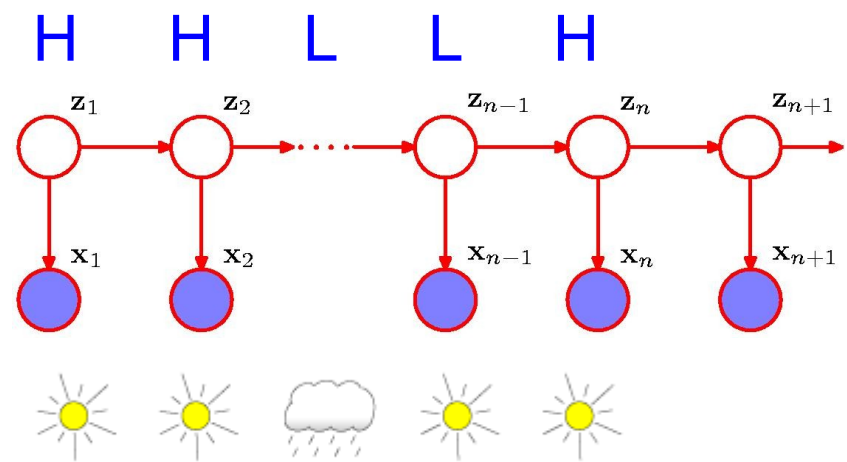
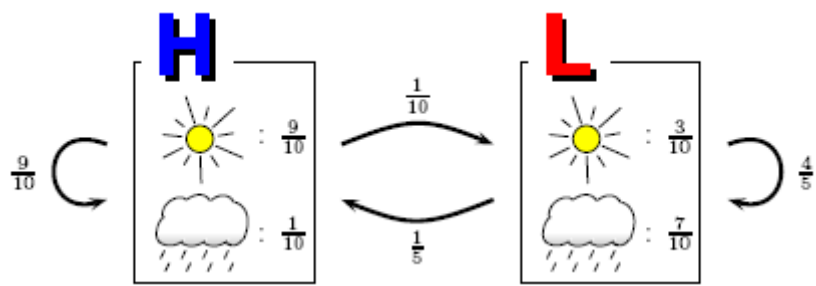
$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

We simply count how many times each outcome of the multinomial variables (a transition or emission) is observed ...

This yield a maximum likelihood estimate (MLE)  $\boldsymbol{\theta}^*$  of  $p(\mathbf{X},\mathbf{Z} \mid \boldsymbol{\theta})$ , which is what we mathematically want ...

**Any problems?** What if e.g. the transition from state  $j$  to  $k$  is *not* observed, then probability  $A_{jk}$  is set to 0. Practical solution: Assume that every transition and emission is seen once (pseudocount) ...

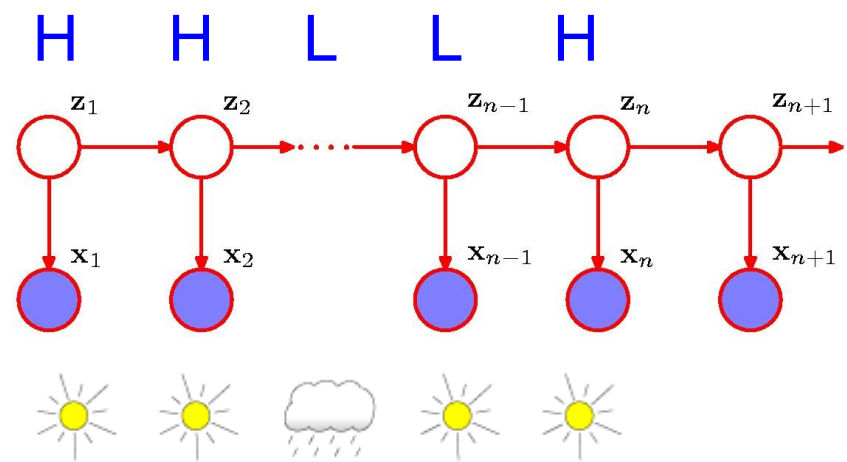
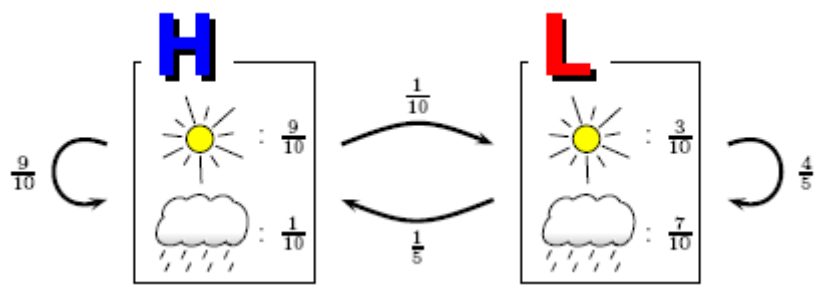
# Example



## Without pseudocounts:

$A_{HH} = 1/2$        $p(\text{sun}|H) = 1$   
 $A_{HL} = 1/2$        $p(\text{rain}|H) = 0$   
 $A_{LH} = 1/2$        $p(\text{sun}|L) = 1/2$   
 $A_{LL} = 1/2$        $p(\text{rain}|L) = 1/2$   
 $\pi_H = 1$   
 $\pi_L = 0$

# Example



## Without pseudocounts:

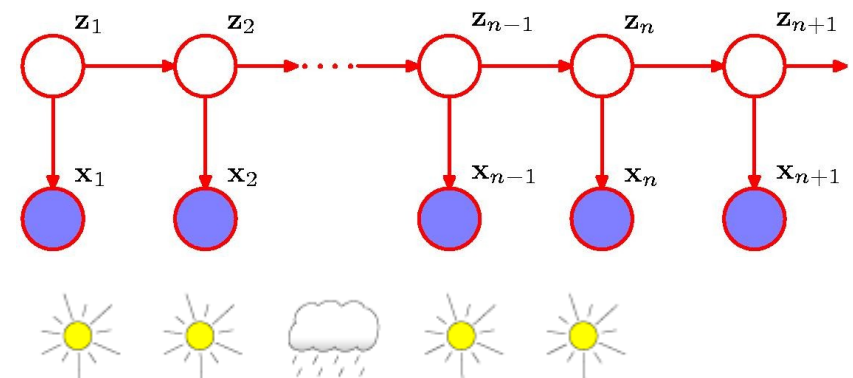
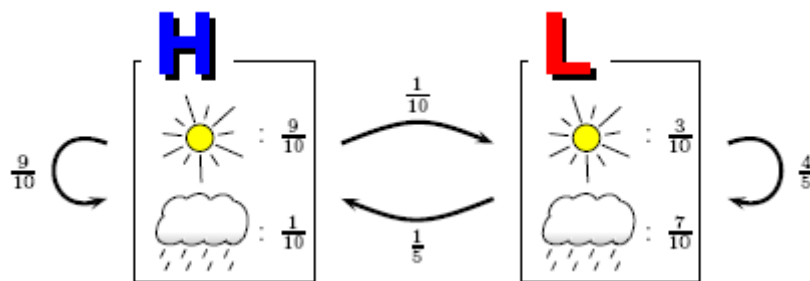
$A_{HH} = 1/2$	$p(\text{sun} H) = 1$
$A_{HL} = 1/2$	$p(\text{rain} H) = 0$
$A_{LH} = 1/2$	$p(\text{sun} L) = 1/2$
$A_{LL} = 1/2$	$p(\text{rain} L) = 1/2$
$\pi_H = 1$	
$\pi_L = 0$	

## With pseudocounts:

$A_{HH} = 2/4$	$p(\text{sun} H) = 4/5$
$A_{HL} = 2/4$	$p(\text{rain} H) = 1/5$
$A_{LH} = 2/4$	$p(\text{sun} L) = 2/4$
$A_{LL} = 2/4$	$p(\text{rain} L) = 2/4$
$\pi_H = 2/3$	
$\pi_L = 1/3$	

# Selecting “the right” parameters

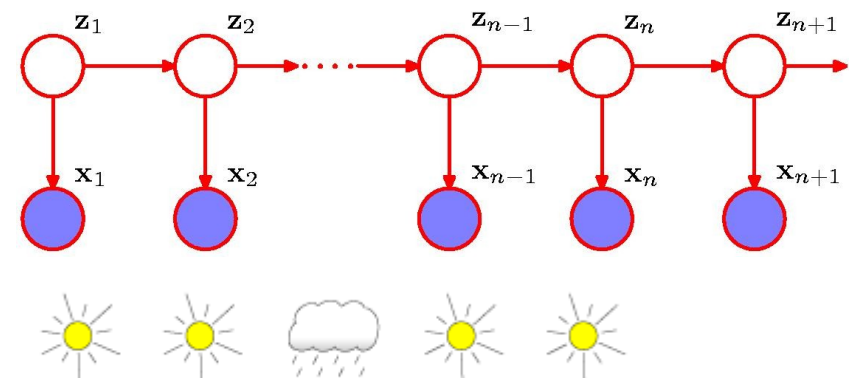
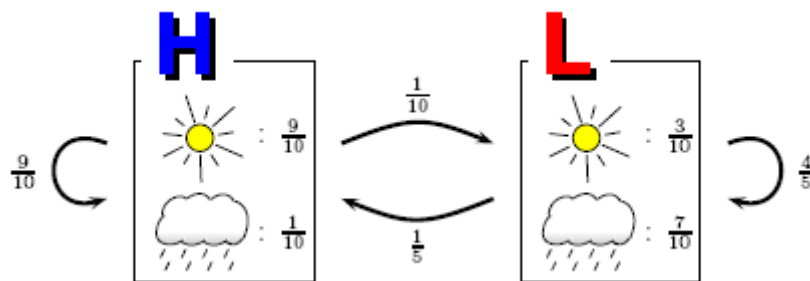
What if only (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  is given, i.e the corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1, \dots, \mathbf{z}_N\}$  are unknown?



How should we set the model parameters, i.e. transitions  $\mathbf{A}$ ,  $\boldsymbol{\pi}$ , and emission probabilities  $\boldsymbol{\Phi}$ , to make the given  $\mathbf{X}$ 's most likely?

# Selecting “the right” parameters

What if only (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  is given, i.e the corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are unknown?



How should we set the model parameters, i.e. transitions  $\mathbf{A}$ ,  $\boldsymbol{\pi}$ , and emission probabilities  $\boldsymbol{\Phi}$ , to make the given  $\mathbf{X}$ 's most likely?

$$\text{Maximize } p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \text{ w.r.t. } \boldsymbol{\theta} \dots$$



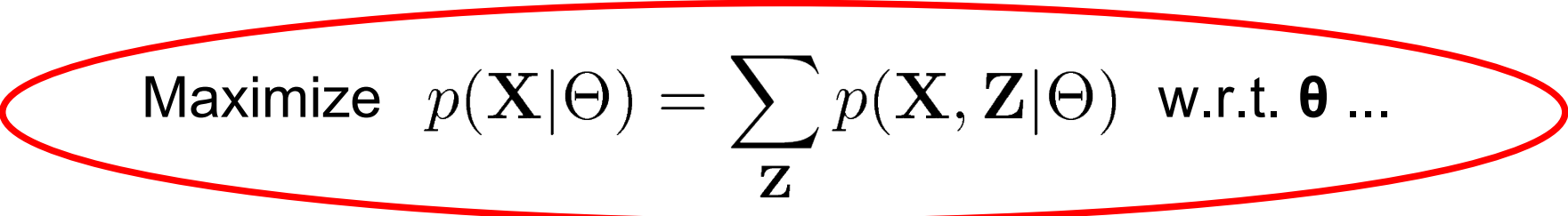
# Selecting “the right” parameters

What if only (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  is given, i.e the corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$  are unknown?

Direct maximization of the likelihood (or *log-likelihood*) is *hard* ...

$$p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$$

How should we set the model parameters, i.e. transitions  $\mathbf{A}$ ,  $\boldsymbol{\pi}$ , and emission probabilities  $\boldsymbol{\Phi}$ , to make the given  $\mathbf{X}$ 's most likely?



Maximize  $p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta)$  w.r.t.  $\boldsymbol{\theta}$  ...

# Practical Solution - Viterbi training

A more “practical” thing to do is **Viterbi Training**:

1. Decide on some initial parameter  $\theta^0$
2. Find the most likely sequence of states  $\mathbf{Z}^*$  explaining  $\mathbf{X}$  using the the Viterbi Algorithm and the current parameters  $\theta^i$
3. Update parameters to  $\theta^{i+1}$  by “counting” (with pseudo counts) according to  $(\mathbf{X}, \mathbf{Z}^*)$ .
4. Repeat 2-3 until  $P(\mathbf{X}, \mathbf{Z}^* | \theta^i)$  is satisfactory (or the Viterbi sequence of states does not change).

# Practical Solution - Viterbi training

A more “practical” thing to do is **Viterbi Training**:

1. Decide on some initial parameter  $\theta^0$
2. Find the most likely sequence of states  $\mathbf{Z}^*$  explaining  $\mathbf{X}$  using the the Viterbi Algorithm and the current parameters  $\theta^i$
3. Update parameters to  $\theta^{i+1}$  by “counting” (with pseudo counts) according to  $(\mathbf{X}, \mathbf{Z}^*)$ .
4. Repeat 2-3 until  $P(\mathbf{X}, \mathbf{Z}^* | \theta^i)$  is satisfactory (or the Viterbi sequence of states does not change).

Finds a (local) maximum of:

$$\text{VIT}_{\mathbf{X}}(\Theta) = \max_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \Theta)$$

The identified parameters  $\theta^*$  is not a MLE of  $p(\mathbf{X} | \theta)$ , but works “ok”

# Summary: Training-by-Counting

**Training-by-Counting:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and the corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ . We want to find a model:

$$\Theta_{\text{TbC}}^* = \arg \max_{\Theta} p(\mathbf{X}, \mathbf{Z} | \Theta) = \arg \max_{\Theta} \log p(\mathbf{X}, \mathbf{Z} | \Theta)$$

This can be done analytically by counting the frequency by which each transition and emission occur in the training data  $(\mathbf{X}, \mathbf{Z})$ .

If only  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is given, then we want to find a model:

$$\Theta_{\mathbf{X}}^* = \arg \max_{\Theta} p(\mathbf{X} | \Theta) = \arg \max_{\Theta} \log p(\mathbf{X} | \Theta)$$

# Summary: Viterbi Training

**Viterbi Training:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . Pick an initial set of parameters  $\boldsymbol{\theta}_{\text{vit}}^0$  and compute the best explanation of  $\mathbf{X}$  under assumption of these parameters using the Viterbi algorithm:

$$\mathbf{Z}_{\text{Vit}}^0 = \arg \max_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^0) = \arg \max_{\mathbf{Z}} \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^0)$$

Compute  $\boldsymbol{\theta}_{\text{vit}}^1$  from  $\boldsymbol{\theta}_{\text{vit}}^0$  and  $\mathbf{Z}_{\text{vit}}^0$  using TbC and iterate:

$$\boldsymbol{\theta}_{\text{Vit}}^1 = \arg \max_{\boldsymbol{\theta}} p(\mathbf{X}, \mathbf{Z}_{\text{Vit}}^0 | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{X}, \mathbf{Z}_{\text{Vit}}^0 | \boldsymbol{\theta})$$

$$\mathbf{Z}_{\text{Vit}}^1 = \arg \max_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^1) = \arg \max_{\mathbf{Z}} \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^1)$$

$p(\mathbf{X} | \boldsymbol{\theta}_{\text{Vit}}^i)$  is usually close to  $p(\mathbf{X} | \boldsymbol{\theta}_{\mathbf{X}}^*)$ , but no guarantees

# Expectation Maximization

**EM Training:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . Pick an initial set of parameters  $\boldsymbol{\theta}_{\text{EM}}^0$  and consider the expectation of  $\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$  over  $\mathbf{Z}$  (given  $\mathbf{X}$  and  $\boldsymbol{\theta}_{\text{EM}}^0$ ) as a function of  $\boldsymbol{\theta}$ :

$$EM_{\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0}(\boldsymbol{\theta}) = E_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0}(\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0) \log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

For HMMs, we can find  $\boldsymbol{\theta}_{\text{EM}}^1$  analytically, and iterate to get  $\boldsymbol{\theta}_{\text{EM}}^i$ :

$$\boldsymbol{\theta}_{\text{EM}}^1 = \arg \max_{\boldsymbol{\theta}} EM_{\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0}(\boldsymbol{\theta})$$

$p(\mathbf{X}|\boldsymbol{\theta}_{\text{EM}}^i)$  converges towards a (local) maximum of  $p(\mathbf{X}|\boldsymbol{\theta})$

# Expectation Maximization

**E-Step:** Define the Q-function:

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

i.e. the expectation of  $\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$  over  $\mathbf{Z}$  (given  $\mathbf{X}$  and  $\boldsymbol{\theta}^{\text{old}}$ ) as a function of  $\boldsymbol{\theta}$

**M-Step:** Maximize  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$  w.r.t.  $\boldsymbol{\theta}$

$$\text{EM}_{\mathbf{X}, \Theta^{\text{old}}}(\Theta) = E_{\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}}[\log p(\mathbf{X}, \mathbf{Z}|\Theta)] = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

$$\Theta^* = \arg \max_{\Theta} \text{EM}_{\mathbf{X}, \Theta^{\text{old}}}(\Theta)$$

When iterated, the likelihood  $p(\mathbf{X}|\boldsymbol{\theta})$  converges to a (local) maximum

# Maximizing the likelihood

Direct maximization of the likelihood (or *log*-likelihood) is *hard* ...

$$p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \qquad p(\mathbf{X}|\Theta) = \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

Assume that we have valid set of parameters  $\Theta^{\text{old}}$ , and that we want to estimate a set  $\Theta$  which yields a better likelihood. We can write:

$$\log p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}|\Theta)$$



# Maximizing the likelihood

Direct maximization of the likelihood (or *log*-likelihood) is *hard* ...

$$p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \qquad p(\mathbf{X}|\Theta) = \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

Assume that we have valid set of parameters  $\Theta^{\text{old}}$ , and that we want to estimate a set  $\Theta$  which yields a better likelihood. We can write:

This sums to 1 ...

$$\log p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}|\Theta)$$

# Maximizing the likelihood

Direct maximization of the likelihood (or *log*-likelihood) is *hard* ...

$$p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \qquad p(\mathbf{X}|\Theta) = \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

Assume that we have valid set of parameters  $\Theta^{\text{old}}$ , and that we want to estimate a set  $\Theta$  which yields a better likelihood. We can write:

$$\begin{aligned} \log p(\mathbf{X}|\Theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}|\Theta) \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) (\log p(\mathbf{X}, \mathbf{Z}|\Theta) - \log p(\mathbf{Z}|\mathbf{X}, \Theta)) \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta) \end{aligned}$$

# Maximizing the likelihood

Direct maximization of the likelihood (or *log*-likelihood) is *hard* ...

$$p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \qquad p(\mathbf{X}|\Theta) = \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

Assume that we have valid set of parameters  $\Theta^{\text{old}}$ , and that we want to estimate a set  $\Theta$  which yields a better likelihood. We can write:

$$\begin{aligned} \log p(\mathbf{X}|\Theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}|\Theta) \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) (\log p(\mathbf{X}, \mathbf{Z}|\Theta) - \log p(\mathbf{Z}|\mathbf{X}, \Theta)) \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta) \end{aligned}$$

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

The expectation (under  $\Theta^{\text{old}}$ ) of the log-likelihood of the complete data (i.e. observations  $\mathbf{X}$  and underlying states  $\mathbf{Z}$ ) as a function of  $\Theta$

# Maximizing the likelihood

Direct maximization of the likelihood (or *log*-likelihood) is *hard* ...

$$p(\mathbf{X}|\Theta) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\Theta) \qquad p(\mathbf{X}|\Theta) = \frac{p(\mathbf{X}, \mathbf{Z}|\Theta)}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

Assume that we have valid set of parameters  $\Theta^{\text{old}}$ , and that we want to estimate a set  $\Theta$  which yields a better likelihood. We can write:

$$\begin{aligned} \log p(\mathbf{X}|\Theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}|\Theta) \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) (\log p(\mathbf{X}, \mathbf{Z}|\Theta) - \log p(\mathbf{Z}|\mathbf{X}, \Theta)) \\ &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta) \\ &= Q(\Theta, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta) \end{aligned}$$

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

# Maximizing the likelihood

Assume that we have valid set of parameters  $\Theta^{\text{old}}$ , and that we want to estimate a set  $\Theta$  which yields a better likelihood. We have:

$$\log p(\mathbf{X}|\Theta) = Q(\Theta, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta)$$

$$\log p(\mathbf{X}|\Theta^{\text{old}}) = Q(\Theta^{\text{old}}, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$$

The increase of the log-likelihood can thus be written as:

$$\log p(\mathbf{X}|\Theta) - \log p(\mathbf{X}|\Theta^{\text{old}}) =$$

$$Q(\Theta, \Theta^{\text{old}}) - Q(\Theta^{\text{old}}, \Theta^{\text{old}}) + \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

# Maximizing the likelihood

Assume that we have valid set of parameters  $\theta^{\text{old}}$ , and that we want to estimate a set  $\theta$  which yields a better likelihood. We have:

$$\log p(\mathbf{X}|\Theta) = Q(\Theta, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta)$$

$$\log p(\mathbf{X}|\Theta^{\text{old}}) = Q(\Theta^{\text{old}}, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$$

The increase of the log-likelihood can thus be written as:

$$\log p(\mathbf{X}|\Theta) - \log p(\mathbf{X}|\Theta^{\text{old}}) =$$

$$Q(\Theta, \Theta^{\text{old}}) - Q(\Theta^{\text{old}}, \Theta^{\text{old}}) + \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log \frac{p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})}{p(\mathbf{Z}|\mathbf{X}, \Theta)}$$

The relative entropy of  $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$  relative to  $p(\mathbf{Z}|\mathbf{X}, \theta)$ , i.e.  $\geq 0$

# Maximizing the likelihood

Assume that we have valid set of parameters  $\theta^{\text{old}}$ , and that we want to estimate a set  $\theta$  which yields a better likelihood. We have:

$$\log p(\mathbf{X}|\Theta) = Q(\Theta, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta)$$

$$\log p(\mathbf{X}|\Theta^{\text{old}}) = Q(\Theta^{\text{old}}, \Theta^{\text{old}}) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}})$$

The increase of the log-likelihood can thus be written as:

$$\log p(\mathbf{X}|\Theta) - \log p(\mathbf{X}|\Theta^{\text{old}}) \geq Q(\Theta, \Theta^{\text{old}}) - Q(\Theta^{\text{old}}, \Theta^{\text{old}})$$

By maximizing the expectation  $Q(\theta, \theta^{\text{old}})$  w.r.t.  $\theta$ , we do not decrease the likelihood, hence name *expectation maximization* ...

# EM for HMMs

**E-Step:** Define the Q-function:

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

i.e. the expectation of  $\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$  over  $\mathbf{Z}$  (given  $\mathbf{X}$  and  $\boldsymbol{\theta}^{\text{old}}$ ) as a function of  $\boldsymbol{\theta}$

**M-Step:** Maximize  $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$  w.r.t.  $\boldsymbol{\theta}$

For HMMs  $Q$  has a closed form and maximization can be performed explicitly. Iterate until no or little increase in likelihood is observed, or some maximum number of iterations is reached ...

When iterated, the likelihood  $p(\mathbf{X}|\boldsymbol{\theta})$  converges to a (local) maximum



# EM for HMMs

- Init:** Pick “suitable” parameters (transition and emission probabilities). Observe that if a parameter is initialized to zero, it remains zero ...
- E-Step:** 1) Run the forward- and backward-algorithms with the current choice of parameters (to get the params of Q-func).
- Stop?:** 2) Compute the likelihood  $p(\mathbf{X}|\boldsymbol{\theta})$ , if sufficient (or another stopping criteria is meet) then stop.
- M-Step:** 3) Compute new parameters using the values stored by the forward- and backward-algorithms. Repeat 1-3.

# EM for HMMs

We want a closed form for  $Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z}|\Theta) &= p(\mathbf{z}_1|\pi) \prod_{n=2}^N p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{z}_n, \phi) \\ &= \prod_{k=1}^K \pi_k^{z_{1k}} \prod_{n=2}^N \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}} \prod_{n=1}^N \prod_{k=1}^K p(\mathbf{x}_n|\phi_k)^{z_{nk}} \end{aligned}$$

# EM for HMMs

We want a closed form for  $Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$

$$\begin{aligned}
 p(\mathbf{X}, \mathbf{Z}|\Theta) &= p(\mathbf{z}_1|\pi) \prod_{n=2}^N p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{z}_n, \phi) \\
 &= \prod_{k=1}^K \pi_k^{z_{1k}} \prod_{n=2}^N \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}} \prod_{n=1}^N \prod_{k=1}^K p(\mathbf{x}_n|\phi_k)^{z_{nk}}
 \end{aligned}$$

$$p(\mathbf{z}_1|\pi) = \prod_{k=1}^K \pi_k^{z_{1k}} \quad p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$$

$$p(\mathbf{x}_n|\mathbf{z}_n, \phi) = \prod_{k=1}^K p(\mathbf{x}_n|\phi_k)^{z_{nk}}$$

# EM for HMMs

We want a closed form for  $Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$

$$\begin{aligned}
 p(\mathbf{X}, \mathbf{Z}|\Theta) &= p(\mathbf{z}_1|\pi) \prod_{n=2}^N p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{z}_n, \phi) \\
 &= \prod_{k=1}^K \pi_k^{z_{1k}} \prod_{n=2}^N \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}} \prod_{n=1}^N \prod_{k=1}^K p(\mathbf{x}_n|\phi_k)^{z_{nk}}
 \end{aligned}$$

Taking the log yields:

$$\log p(\mathbf{X}, \mathbf{Z}|\Theta) = \sum_{k=1}^K z_{1k} \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K z_{n-1,j} z_{nk} \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K z_{nk} \log p(\mathbf{x}_n|\phi_k)$$

# EM for HMMs

We want a closed form for  $Q(\Theta, \Theta^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \Theta^{\text{old}}) \log p(\mathbf{X}, \mathbf{Z}|\Theta)$

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z}|\Theta) &= p(\mathbf{z}_1|\pi) \prod_{n=2}^N p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}) \prod_{n=1}^N p(\mathbf{x}_n|\mathbf{z}_n, \phi) \\ &= \prod_{k=1}^K \pi_k^{z_{1k}} \prod_{n=2}^N \prod_{j=1}^K \prod_{k=1}^K A_{jk}^{z_{n-1,j} z_{nk}} \prod_{n=1}^N \prod_{k=1}^K p(\mathbf{x}_n|\phi_k)^{z_{nk}} \end{aligned}$$

Taking the log yields:

$$\log p(\mathbf{X}, \mathbf{Z}|\Theta) = \sum_{k=1}^K z_{1k} \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K z_{n-1,j} z_{nk} \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K z_{nk} \log p(\mathbf{x}_n|\phi_k)$$

Taking the expectation (under  $\Theta^{\text{old}}$  and  $\mathbf{X}$ ) over  $\mathbf{Z}$  yields  $Q(\Theta, \Theta^{\text{old}})$ , i.e:

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K E(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K E(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K E(z_{nk}) \log p(\mathbf{x}_n|\phi_k)$$

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K E(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K E(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K E(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**E-Step:** To calculate  $Q$ , we must compute the expectations  $E(z_{1k})$ ,  $E(z_{nk})$ , and  $E(z_{n-1,j} z_{nk})$ . Consider the probabilities:

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \Theta^{\text{old}})$$

A  $K$ -vector where entry  $k$  is the prob  $\gamma(z_{nk})$  of being in state  $k$  in the  $n$ 'th step ...

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \Theta^{\text{old}})$$

A  $K \times K$ -table where entry  $(j, k)$  is the prob  $\xi(z_{n-1,j} z_{nk})$  of being in state  $j$  and  $k$  in the  $(n-1)$ 'th and  $n$ 'th step ...

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K E(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K E(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K E(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**E-Step:** To calculate  $Q$ , we must compute the expectations  $E(z_{1k})$ ,  $E(z_{nk})$ , and  $E(z_{n-1,j} z_{nk})$ . Consider the probabilities:

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \Theta^{\text{old}})$$

A  $K$ -vector where entry  $k$  is the prob  $\gamma(z_{nk})$  of being in state  $k$  in the  $n$ 'th step ...

binary variables

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \Theta^{\text{old}})$$

A  $K \times K$ -table where entry  $(j,k)$  is the prob  $\xi(z_{n-1,j} z_{nk})$  of being in state  $j$  and  $k$  in the  $(n-1)$ 'th and  $n$ 'th step ...

**Fact:** The expectation of a binary variable  $z$  is just  $p(z=1)$  ...

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K E(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K E(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K E(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**E-Step:** To calculate Q, we need  $E(z_{nk})$ , and  $E(z_{n-1,j} z_{nk})$ . Cor

$$E(z_{nk}) = \gamma(z_{nk})$$
$$E(z_{n-1,j} z_{nk}) = \xi(z_{n-1,j} z_{nk})$$

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \Theta^{\text{old}})$$

prob  $\gamma(\mathbf{z}_{nk})$  of being in state k in the n'th step ...

binary variables

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \Theta^{\text{old}})$$

A  $K \times K$ -table where  $\xi(z_{n-1,j}, z_{nk})$  is the prob  $\xi(z_{n-1,j} z_{nk})$  of being in state j and k in the (n-1)'th and n'th step ...

**Fact:** The expectation of a binary variable z is just  $p(z=1)$  ...



# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**M-Step:** If we assume discrete observables  $x_i$ , then maximizing the above w.r.t.  $\Theta$ , i.e.  $A$ ,  $\pi$ , and  $\Phi$ , yields:

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**M-Step:** If we assume discrete observables  $x_i$ , then maximizing the above w.r.t.  $\theta$ , i.e.  $A$ ,  $\pi$ , and  $\Phi$ , yields:

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j} z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j} z_{nl})}$$

Expected number of transitions from state  $j$  to state  $k$

Expected number of transitions from state  $j$  to any state

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})}$$

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**M-Step:** If we assume discrete observables  $x_i$ , then maximizing the above w.r.t.  $\Theta$ , i.e.  $A$ ,  $\pi$ , and  $\Phi$ , yields:

$$\phi_{ik} = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\sum_{n=1}^N \gamma(z_{nk})} =$$

Expected number of times  
symbol  $i$  is emitted from state  $k$

Expected number of times a  
symbol is emitted from state  $k$

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**M-Step:** If we assume discrete observables  $x_i$ , then maximizing the above w.r.t.  $\Theta$ , i.e.  $A$ ,  $\pi$ , and  $\Phi$ , yields:

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} \quad \pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} \quad \phi_{ik} = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\sum_{n=1}^N \gamma(z_{nk})}$$

# EM for HMMs

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{k=1}^K \gamma(z_{1k}) \log \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j} z_{nk}) \log A_{jk} + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \log p(\mathbf{x}_n | \phi_k)$$

**M-Step:** If we assume discrete observables  $x_i$ , then maximizing the above w.r.t.  $\Theta$ , i.e.  $A$ ,  $\pi$ , and  $\Phi$ , yields:

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} \quad \pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} \quad \phi_{ik} = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\sum_{n=1}^N \gamma(z_{nk})}$$

Compare this to the formulas when  $\mathbf{X}$  and  $\mathbf{Z}$  where given:

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

# Computing $\gamma$ and $\xi$

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}) \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X})}\end{aligned}$$

$$\begin{aligned}\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}) \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_{n-1}) \beta(\mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n)}{p(\mathbf{X})}\end{aligned}$$

Can be computed efficiently using the forward- and backward-algorithm

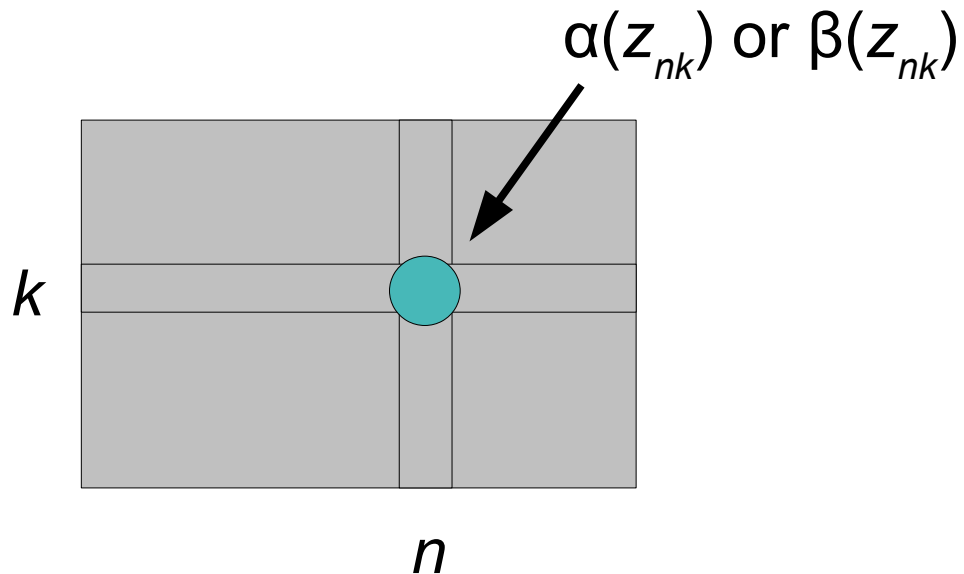
# Computing the new parameters

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} = \frac{\alpha(z_{1k})\beta(z_{1k})/p(\mathbf{X})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})/p(\mathbf{X})} = \frac{\alpha(z_{1k})\beta(z_{1k})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})}$$

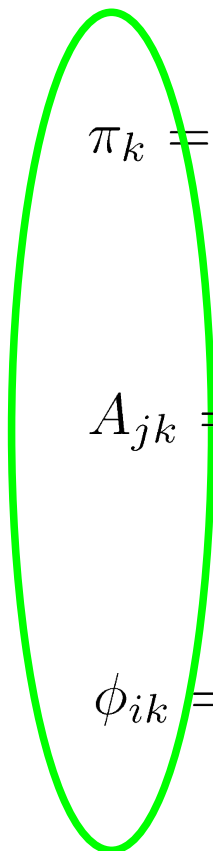
$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} = \frac{\sum_{n=2}^N \alpha(z_{n-1,j})\beta(z_{nk})p(\mathbf{x}_n|\phi_k)A_{jk}}{\sum_{l=1}^K \sum_{n=2}^N \alpha(z_{n-1,j})\beta(z_{nl})p(\mathbf{x}_n|\phi_l)A_{jl}}$$

$$\phi_{ik} = \frac{\sum_{n=1}^N \alpha(z_{nk})\beta(z_{nk})x_{ni}}{\sum_{n=1}^N \alpha(z_{nk})\beta(z_{nk})}$$

$$\gamma(\mathbf{z}_n) = \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})}$$
$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = \frac{\alpha(\mathbf{z}_{n-1})\beta(\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)}{p(\mathbf{X})}$$



# Computing the new parameters



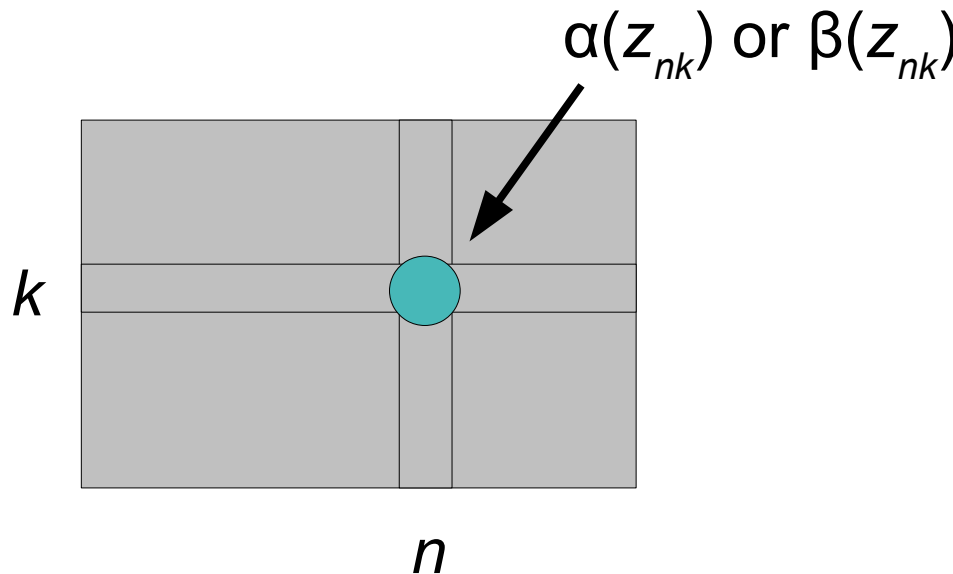
The old parameters

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} = \frac{\alpha(z_{1k})\beta(z_{1k})/p(\mathbf{X})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})/p(\mathbf{X})} = \frac{\alpha(z_{1k})\beta(z_{1k})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})}$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} = \frac{\sum_{n=2}^N \alpha(z_{n-1,j})\beta(z_{nk})p(\mathbf{x}_n|\phi_k)A_{jk}}{\sum_{l=1}^K \sum_{n=2}^N \alpha(z_{n-1,j})\beta(z_{nl})p(\mathbf{x}_n|\phi_l)A_{jl}}$$

$$\phi_{ik} = \frac{\sum_{n=1}^N \alpha(z_{nk})\beta(z_{nk})x_{ni}}{\sum_{n=1}^N \alpha(z_{nk})\beta(z_{nk})}$$

The new parameters





# EM for HMMs - Summary

- Init:** Pick “suitable” parameters (transition and emission probabilities). Observe that if a parameter is initialized to zero, it remains zero ...
- E-Step:** 1) Run the forward- and backward-algorithms with the current choice of parameters (to get the params of Q-func).
- Stop?:** 2) Compute the likelihood  $p(\mathbf{X}|\boldsymbol{\theta})$ , if sufficient (or another stopping criteria is met) then stop.
- M-Step:** 3) Compute new parameters using the values stored by the forward- and backward-algorithms. Repeat 1-3.

Running time per iteration:

$O(K^2N + KK + K^2NK + KDN)$ , where  $D$  is number of observable symbols

By using memorization in 3), we can improve it to  $O(K^2N + KDN)$

# Using the scaled values in EM

$$\begin{aligned}
 \gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}) \\
 &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} \\
 &= \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X})} \\
 &= \hat{\alpha}(\mathbf{z}_n) \hat{\beta}(\mathbf{z}_n)
 \end{aligned}$$

$$\begin{aligned}
 \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}) \\
 &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} \\
 &= \frac{\alpha(\mathbf{z}_{n-1}) \beta(\mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n)}{p(\mathbf{X})} \\
 &= \hat{\alpha}(\mathbf{z}_{n-1}) \hat{\beta}(\mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) / c_n
 \end{aligned}$$

Can be computed using the modified forward- and backward-algorithm

# Using the scaled val

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n|\mathbf{X}) \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n)p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N|\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \hat{\alpha}(\mathbf{z}_n)\hat{\beta}(\mathbf{z}_n)\end{aligned}$$

$$\begin{aligned}p(\mathbf{X}) &= \prod_{n=1}^N c_n \\ \alpha(\mathbf{z}_n) &= \left(\prod_{m=1}^n c_m\right) \hat{\alpha}(\mathbf{z}_n) \\ \beta(\mathbf{z}_n) &= \left(\prod_{m=n+1}^N c_m\right) \hat{\beta}(\mathbf{z}_n)\end{aligned}$$

$$\begin{aligned}\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n|\mathbf{X}) \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1})p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N|\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\alpha(\mathbf{z}_{n-1})\beta(\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \hat{\alpha}(\mathbf{z}_{n-1})\hat{\beta}(\mathbf{z}_n)p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)/c_n\end{aligned}$$

Error in book

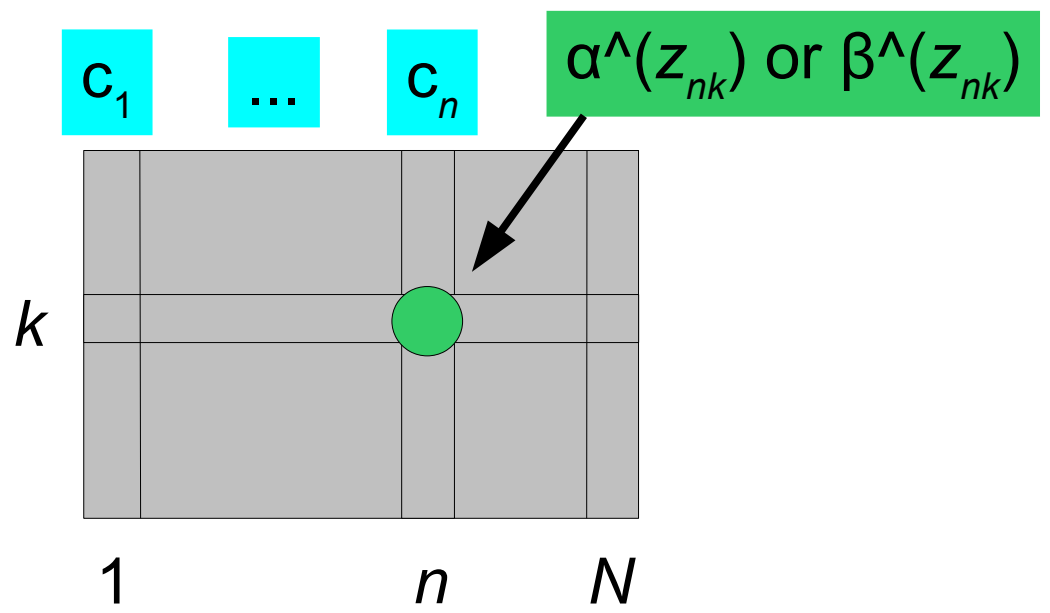
Can be computed using the modified forward- and backward-algorithm

# Computing the new parameters

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} = \frac{\hat{\alpha}(z_{1k})\hat{\beta}(z_{1k})}{\sum_{j=1}^K \hat{\alpha}(z_{1j})\hat{\beta}(z_{1j})}$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} = \frac{\sum_{n=2}^N \hat{\alpha}(z_{n-1,j})\hat{\beta}(z_{nk})p(\mathbf{x}_n|\phi_k)A_{jk}/c_n}{\sum_{l=1}^K \sum_{n=2}^N \hat{\alpha}(z_{n-1,j})\hat{\beta}(z_{nl})p(\mathbf{x}_n|\phi_l)A_{jl}/c_n}$$

$$\phi_{ik} = \frac{\sum_{n=1}^N \hat{\alpha}(z_{nk})\hat{\beta}(z_{nk})x_{ni}}{\sum_{n=1}^N \hat{\alpha}(z_{nk})\hat{\beta}(z_{nk})}$$



# Summary

- Selecting parameters by **counting** to reflect a set of  $(\mathbf{X}, \mathbf{Z})$ 's, i.e. if full information about observables and corresponding latent values is given.
- Selecting parameters by **Viterbi Training** or **Expectation Maximization** to reflect a set of  $\mathbf{X}$ 's, i.e. if only information about observables is given.

# Summary

- Selecting parameters by **counting** to reflect a set of  $(\mathbf{X}, \mathbf{Z})$ 's, i.e. if full information about observables and corresponding latent values is given.
- Selecting parameters by **Viterbi Training** or **Expectation Maximization** to reflect a set of  $\mathbf{X}$ 's, i.e. if only information about observables is given.

How to deal with multiple “training sequences”?

# When multiple $(\mathbf{X}, \mathbf{Z})$ 's are given ...

Assume that (several) sequences of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1, \dots, \mathbf{z}_N\}$  are given ...

$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

... just sum each nominator and denominator over all  $(\mathbf{X}, \mathbf{Z})$ 's, i.e. we divide total counts ...

$$A_{jk} = \frac{\sum_{(\mathbf{X}, \mathbf{Z})} \sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{(\mathbf{X}, \mathbf{Z})} \sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{\sum_{(\mathbf{X}, \mathbf{Z})} z_{1k}}{\sum_{(\mathbf{X}, \mathbf{Z})} \sum_{j=1}^K z_{1j}}$$

$$\phi_{ik} = \frac{\sum_{(\mathbf{X}, \mathbf{Z})} \sum_{n=1}^N z_{nk} x_{ni}}{\sum_{(\mathbf{X}, \mathbf{Z})} \sum_{n=1}^N z_{nk}}$$

# When multiple $\mathbf{X}$ 's are given ...

Assume that a set sequences of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  is given

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} \quad \pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} \quad \phi_{ik} = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\sum_{n=1}^N \gamma(z_{nk})}$$

... just sum each nominator and denominator over all  $\mathbf{X}$ 's, i.e. we divide total expectation, and we must run the forward- and backward algorithms for each training sequence  $\mathbf{X}$  ...

$$A_{jk} = \frac{\sum_{\mathbf{X}} \sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{\mathbf{X}} \sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} \quad \pi_k = \frac{\sum_{\mathbf{X}} \gamma(z_{1k})}{\sum_{\mathbf{X}} \sum_{j=1}^K \gamma(z_{1j})}$$

$$\phi_{ik} = \frac{\sum_{\mathbf{X}} \sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\sum_{\mathbf{X}} \sum_{n=1}^N \gamma(z_{nk})}$$



# Summary: Training-by-Counting

**Training-by-Counting:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and the corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$ . We want to find a model:

$$\Theta_{\text{TbC}}^* = \arg \max_{\Theta} p(\mathbf{X}, \mathbf{Z}|\Theta) = \arg \max_{\Theta} \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

This can be done analytically by counting the frequency by which each transition and emission occur in the training data  $(\mathbf{X}, \mathbf{Z})$ .

# Summary: Training-by-Counting

**Training-by-Counting:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_N\}$  and the corresponding latent states  $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_N\}$ . We want to find a model:

$$\Theta_{\text{TbC}}^* = \arg \max_{\Theta} p(\mathbf{X}, \mathbf{Z}|\Theta) = \arg \max_{\Theta} \log p(\mathbf{X}, \mathbf{Z}|\Theta)$$

This can be done analytically by counting the frequency by which each transition and emission occur in the training data  $(\mathbf{X}, \mathbf{Z})$ .

If only  $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_n\}$  is given, then we want to find a model:

$$\Theta_{\mathbf{X}}^* = \arg \max_{\Theta} p(\mathbf{X}|\Theta) = \arg \max_{\Theta} \log p(\mathbf{X}|\Theta)$$

Finding  $\Theta_{\mathbf{X}}^*$  is hard. We have seen two approaches.

# Summary: Viterbi Training

**Viterbi Training:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . Pick an initial set of parameters  $\boldsymbol{\theta}_{\text{vit}}^0$  and compute the best explanation of  $\mathbf{X}$  under assumption of these parameters using the Viterbi algorithm:

$$\mathbf{Z}_{\text{Vit}}^0 = \arg \max_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^0) = \arg \max_{\mathbf{Z}} \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^0)$$

Compute  $\boldsymbol{\theta}_{\text{vit}}^1$  from  $\boldsymbol{\theta}_{\text{vit}}^0$  and  $\mathbf{Z}_{\text{Vit}}^0$  using TbC and iterate:

$$\boldsymbol{\theta}_{\text{Vit}}^1 = \arg \max_{\boldsymbol{\theta}} p(\mathbf{X}, \mathbf{Z}_{\text{Vit}}^0 | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{X}, \mathbf{Z}_{\text{Vit}}^0 | \boldsymbol{\theta})$$

$$\mathbf{Z}_{\text{Vit}}^1 = \arg \max_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^1) = \arg \max_{\mathbf{Z}} \log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}_{\text{vit}}^1)$$

$p(\mathbf{X} | \boldsymbol{\theta}_{\text{Vit}}^i)$  is usually close to  $p(\mathbf{X} | \boldsymbol{\theta}_{\mathbf{X}}^*)$ , but no guarantees

# Summary: Expectation Maximization

**EM Training:** We are given a sequence of observations  $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ . Pick an initial set of parameters  $\boldsymbol{\theta}_{\text{EM}}^0$  and consider the expectation of  $\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta})$  over  $\mathbf{Z}$  (given  $\mathbf{X}$  and  $\boldsymbol{\theta}_{\text{EM}}^0$ ) as a function of  $\boldsymbol{\theta}$ :

$$EM_{\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0}(\boldsymbol{\theta}) = E_{\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0}(\log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0) \log p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

For HMMs, we can find  $\boldsymbol{\theta}_{\text{EM}}^1$  analytically, and iterate to get  $\boldsymbol{\theta}_{\text{EM}}^i$  :

$$\boldsymbol{\theta}_{\text{EM}}^1 = \arg \max_{\boldsymbol{\theta}} EM_{\mathbf{X}, \boldsymbol{\theta}_{\text{EM}}^0}(\boldsymbol{\theta})$$

$p(\mathbf{X}|\boldsymbol{\theta}_{\text{EM}}^i)$  converges towards a (local) maximum of  $p(\mathbf{X}|\boldsymbol{\theta})$