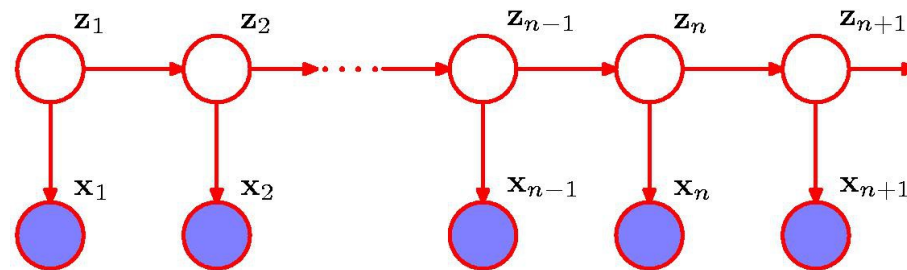


Hidden Markov Models

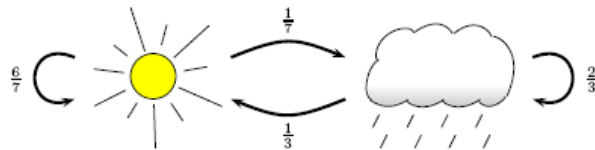
Selecting model parameters or “training”



Hidden Markov Models

Motivation: The n^{th} observation in a chain of observations is influenced by a corresponding latent variable ...

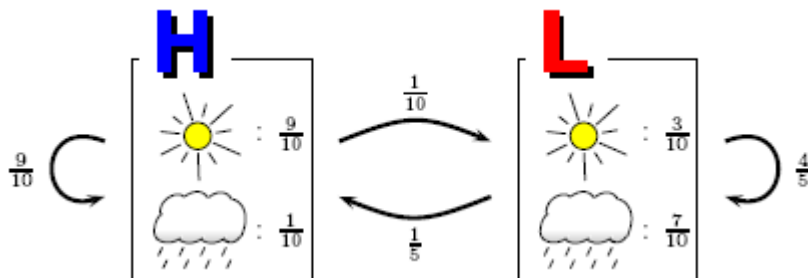
Markov Model



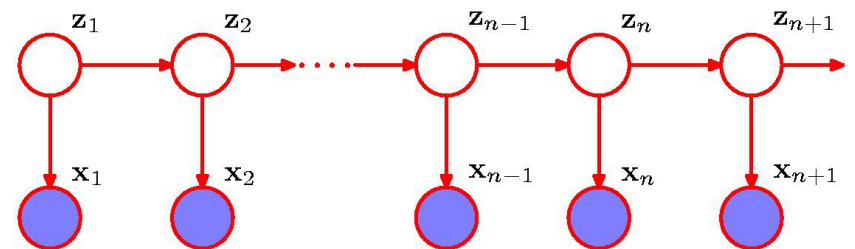
Observations



Hidden Markov Model



Latent values



If the latent states are discrete and form a Markov chain, then it is a **hidden Markov model (HMM)**

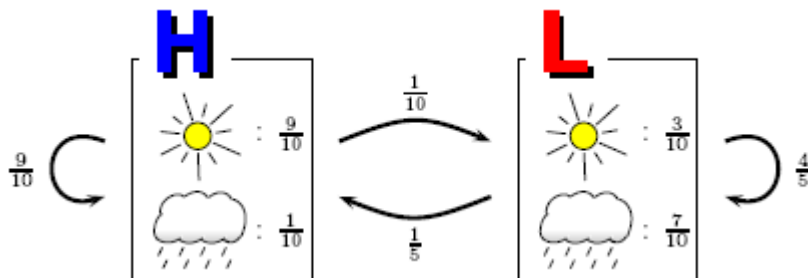
Hidden Markov Models

Motivation: The n^{th} observation in a chain of observations is

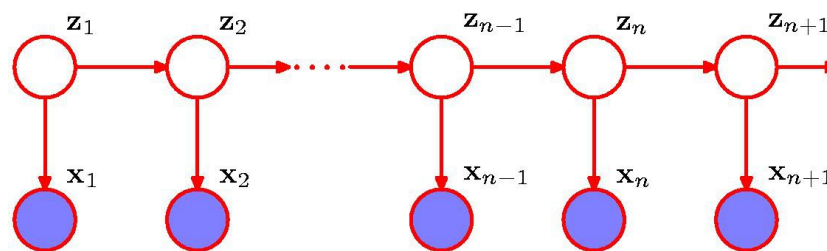
The joint distribution of observables \mathbf{X} and latent values \mathbf{Z} :

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) = p(\mathbf{z}_1) \left[\prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n)$$

Hidden Markov Model



Latent values



If the latent states are discrete and form a Markov chain, then it is a **hidden Markov model (HMM)**

Hidden Markov Models

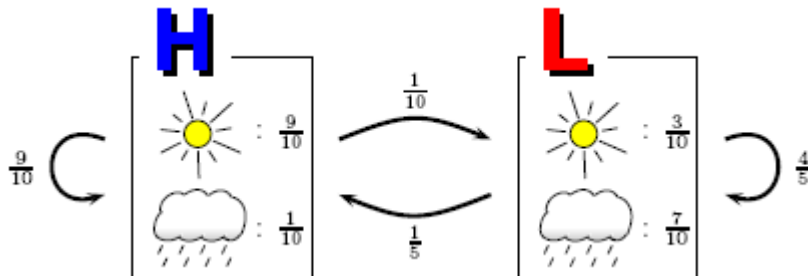
Transition probabilities A and π

Emission probabilities Φ

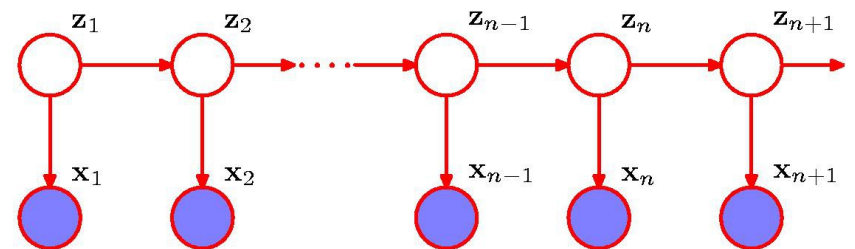
in a chain of observables X and latent values Z :

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) = p(\mathbf{z}_1) \left[\prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n)$$

Hidden Markov Model



Latent values

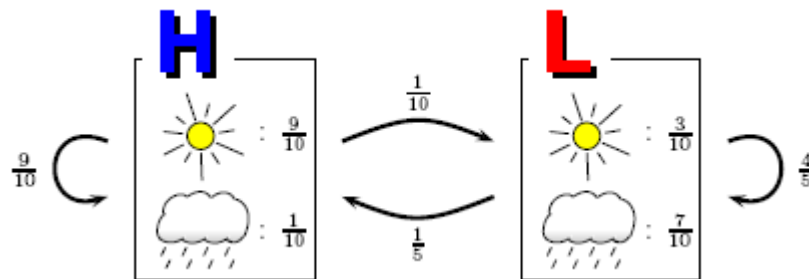


If the latent states are discrete and form a Markov chain, then it is a **hidden Markov model (HMM)**

HMMs as a generative model

A HMM *generates a sequence of observables* by moving from latent state to latent state according to the transition probabilities and *emitting an observable* (from a discrete set of observables, i.e. a finite alphabet) from each latent state visited *according to the emission probabilities* of the state ...

Model M :



A *run* follows a sequence of states:

H H L L H

And *emits* a sequence of symbols:



For a HMM that generates finite strings (e.g. a HMM with an end-state), the language $L = \{\mathbf{X} \mid p(\mathbf{X}) > 0\}$ is regular ...

What we know

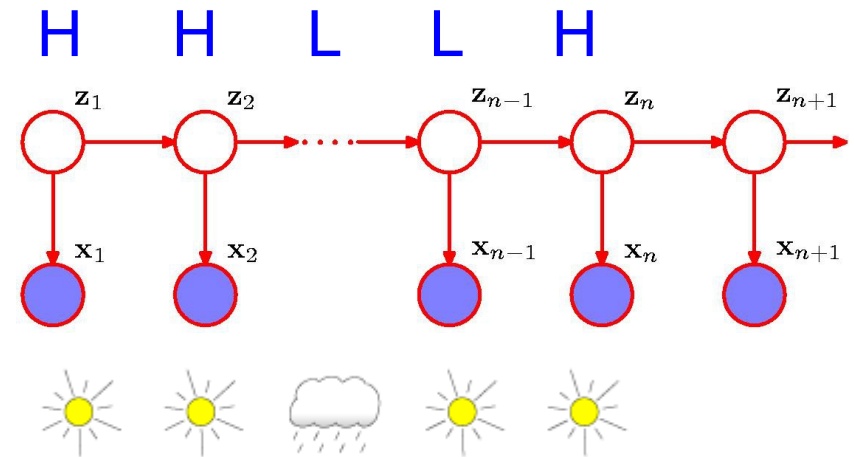
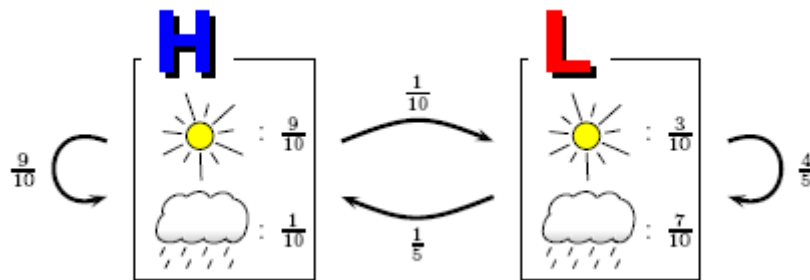
- Introduced hidden Markov models (**HMMs**)
- The **forward- and backward-algorithms** for determining the likelihood $p(\mathbf{X})$ of a sequence of observations, and predicting the next observation in a sequence of observations.
- The **Viterbi-algorithm** for finding the most likely underlying explanation (sequence of latent states) of a sequence of observation
- How to implement them using log-space and scaling.

Today

- 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

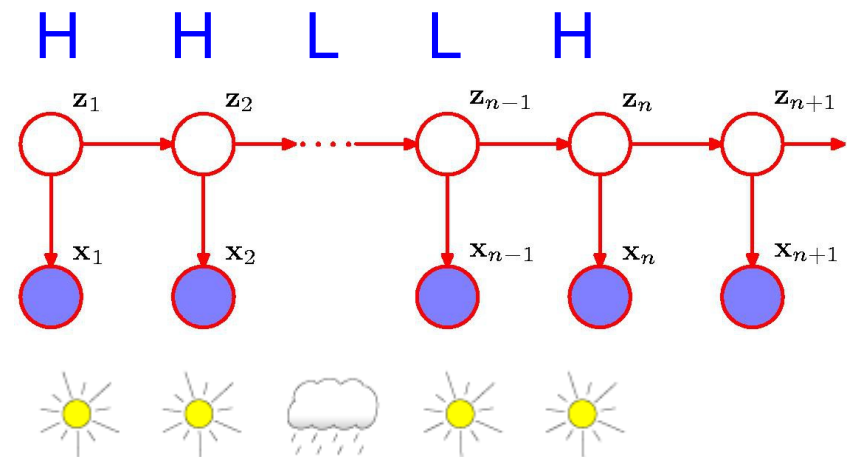
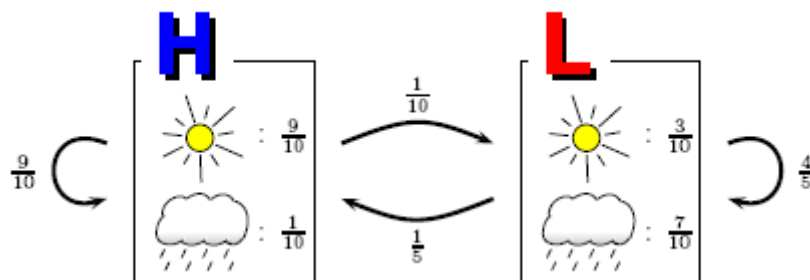
Assume that (several) sequences of observations $\mathbf{X}=\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ and corresponding latent states $\mathbf{Z}=\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ are given ...



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

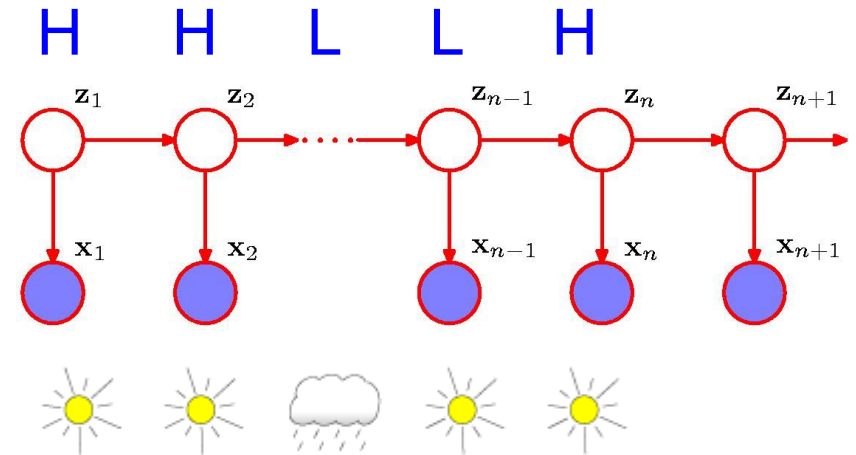
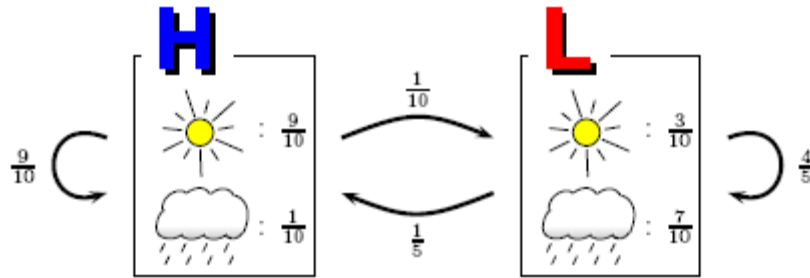
Assume that (several) sequences of observations $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_n\}$ and corresponding latent states $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_n\}$ are given ...



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



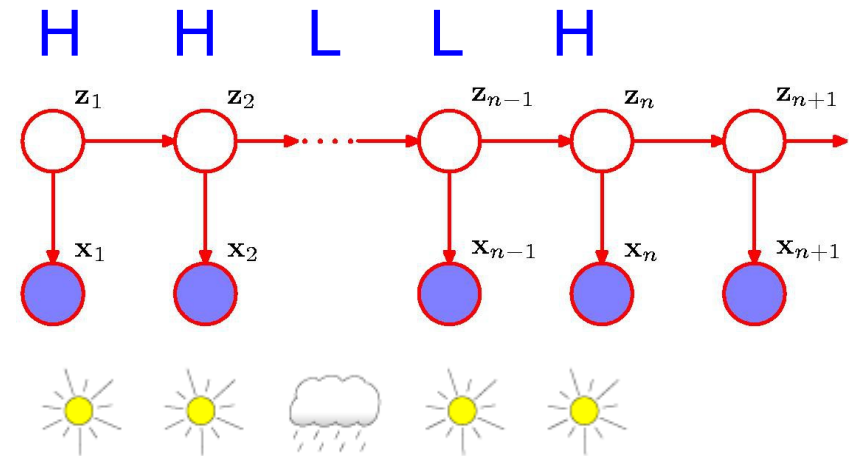
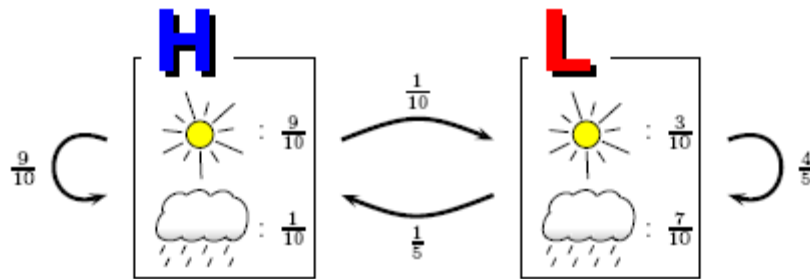
A_{jk} is the probability of a transition from state j to state k , and π_k is the probability of starting in state k ...

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

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 π_k is the probability of starting in state k ...

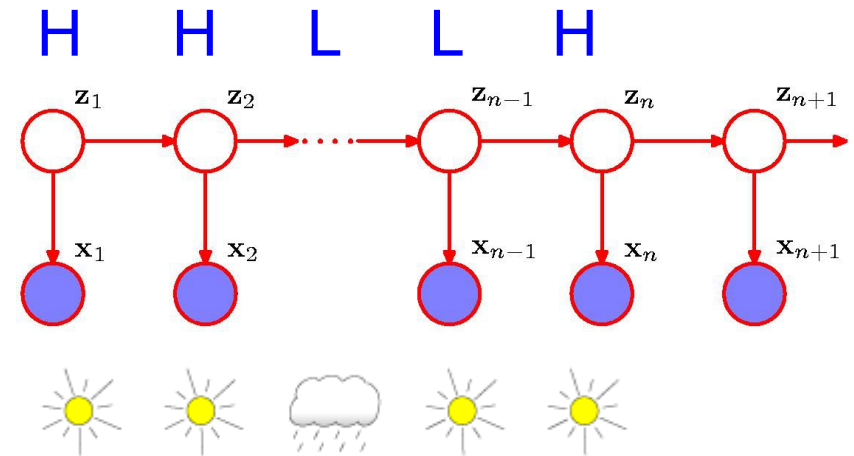
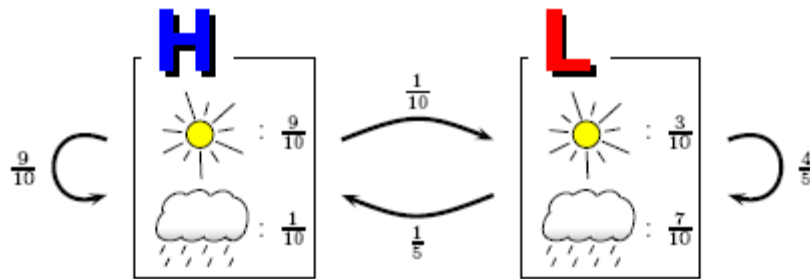
$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} =$$

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

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 ϕ_{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 ...

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

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 yields 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 yields 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 ...

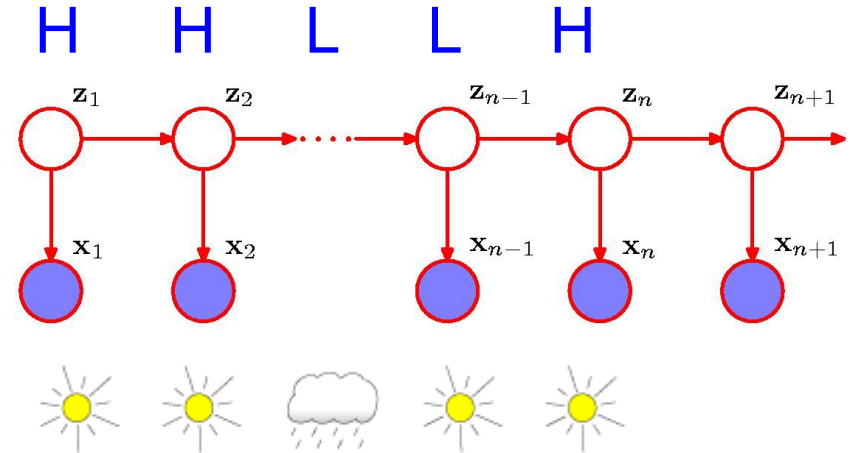
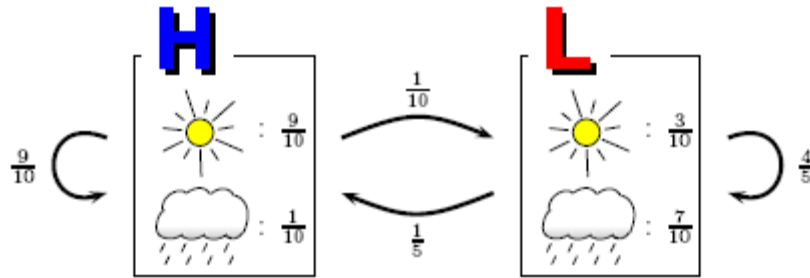
$$A_{jk} = \frac{\sum_{n=2}^N z_{n-1,j} z_{nk}}{\sum_{l=1}^K \sum_{n=2}^N z_{n-1,j} z_{nl}} \quad \pi_k = \frac{z_{1k}}{\sum_{j=1}^K z_{1j}} \quad \phi_{ik} = \frac{\sum_{n=1}^N z_{nk} x_{ni}}{\sum_{n=1}^N z_{nk}}$$

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

This yields 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) ...

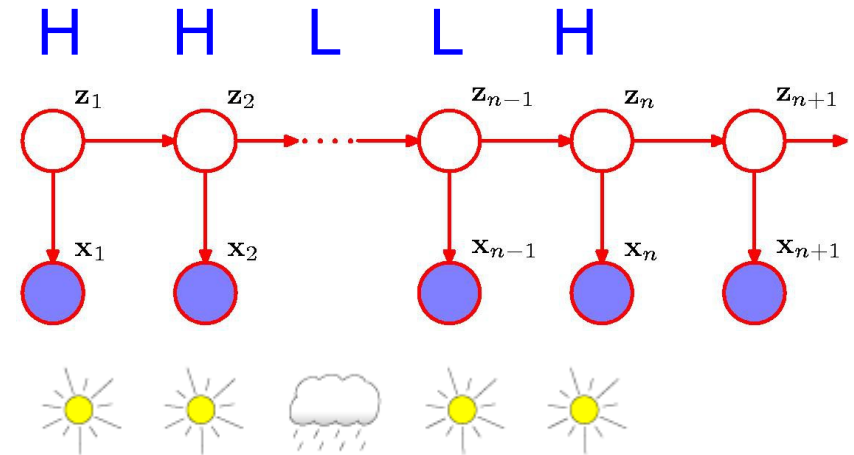
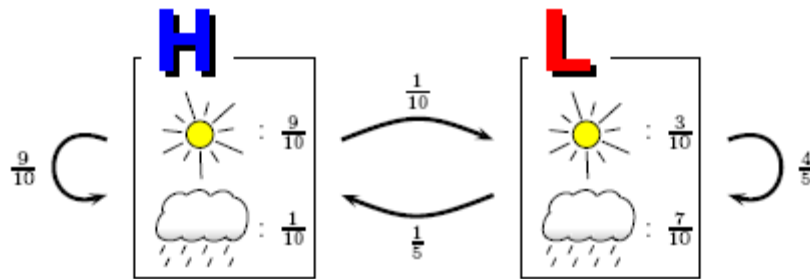
Example



Without pseudocounts:

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

Example



Without pseudocounts:

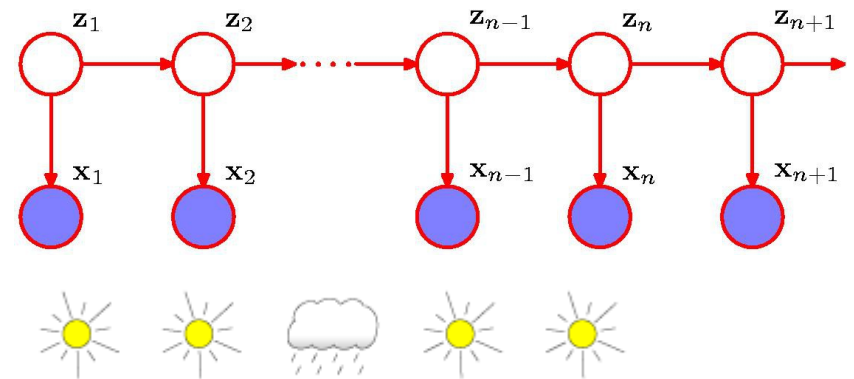
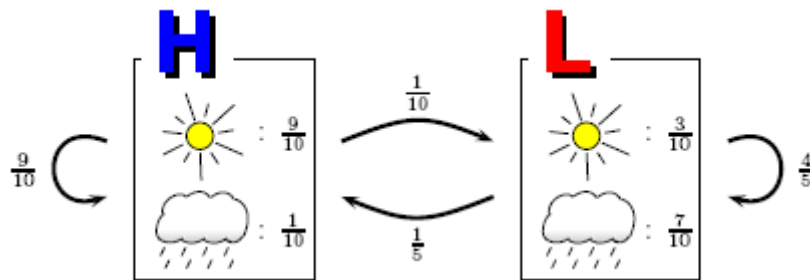
$$\begin{aligned}
 A_{HH} &= 1/2 & p(\text{sun}|H) &= 1 \\
 A_{HL} &= 1/2 & p(\text{rain}|H) &= 0 \\
 A_{LH} &= 1/2 & p(\text{sun}|L) &= 1/2 \\
 A_{LL} &= 1/2 & p(\text{rain}|L) &= 1/2 \\
 \pi_H &= 1 \\
 \pi_L &= 0
 \end{aligned}$$

With pseudocounts:

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

Selecting “the right” parameters

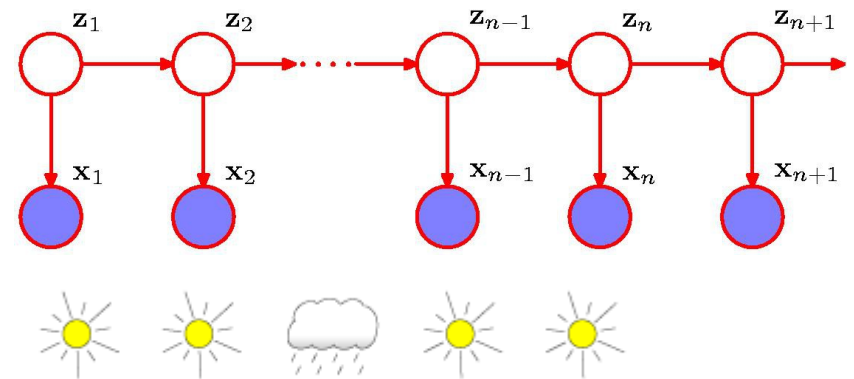
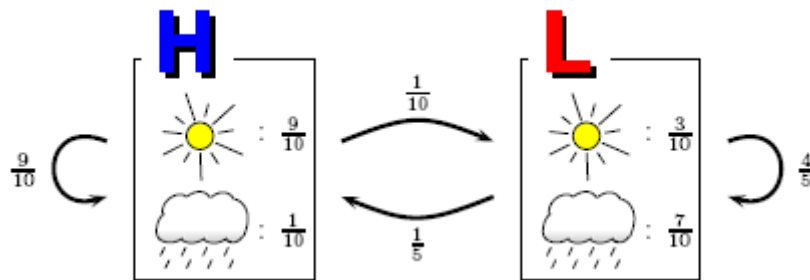
What if only (several) sequences of observations $\mathbf{X}=\{\mathbf{x}_1,\dots,\mathbf{x}_n\}$ is given, i.e the corresponding latent states $\mathbf{Z}=\{\mathbf{z}_1,\dots,\mathbf{z}_n\}$ are unknown?



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?

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

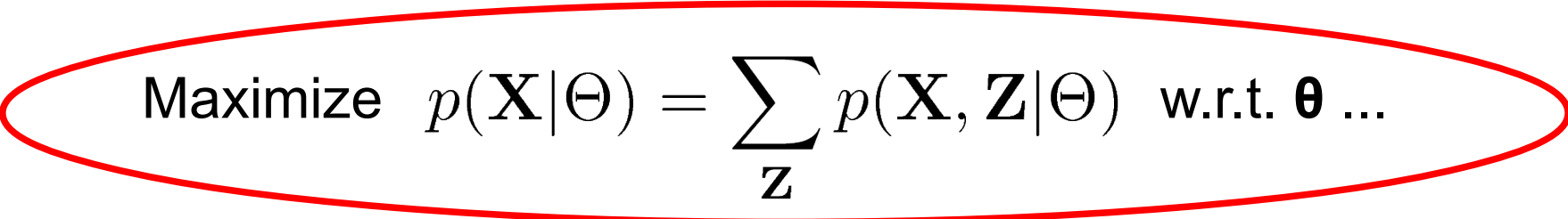
Selecting “the right” parameters

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* ...

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

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?


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

Viterbi training

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

1. Decide on some initial parameter θ^0
2. Find the most likely sequence of states \mathbf{Z}^* explaining \mathbf{X} using the Viterbi Algorithm and the current parameters θ^i
3. Update parameters to θ^{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.

Viterbi training

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

1. Decide on some initial parameter θ^0
2. Find the most likely sequence of states \mathbf{Z}^* explaining \mathbf{X} using the Viterbi Algorithm and the current parameters θ^i
3. Update parameters to θ^{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.

Finds a (local) maximum of:

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

Not a MLE (because right-hand side isn't a likelihood), but works ok

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 the log-likelihood of the complete data (i.e. observations \mathbf{X} and underlying states \mathbf{Z}) as a function of $\boldsymbol{\theta}$

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

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

$$\Theta^* = \arg \max_{\Theta} f_{\text{EM}}(\Theta)$$

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 Θ^{old} , and that we want to estimate a set Θ 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 Θ^{old} , and that we want to estimate a set Θ which yields a better likelihood. We can write:

This sums to 1 ...

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

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 Θ^{old} , and that we want to estimate a set Θ 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 of the log-likelihood of the complete data (i.e. observations \mathbf{X} and underlying states \mathbf{Z}) as a function of Θ

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 Θ^{old} , and that we want to estimate a set Θ 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 Θ^{old} , and that we want to estimate a set Θ 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 θ^{old} , and that we want to estimate a set θ 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. ≥ 0

Maximizing the likelihood

Assume that we have valid set of parameters Θ^{old} , and that we want to estimate a set Θ 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. Θ , we increase 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 the log-likelihood of the complete data (i.e. observations \mathbf{X} and underlying states \mathbf{Z}) as a function of Θ

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

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 $p(\mathbf{X}|\Theta)$ 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 over all \mathbf{Z} 's 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

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

$$E(z_{nk}) = \gamma(z_{nk}) = \sum_{\mathbf{z}} \gamma(\mathbf{z}_n) z_{nk}$$

$$E(z_{n-1,j} z_{nk}) = \xi(z_{n-1,j} z_{nk}) = \sum_{\mathbf{z}} \xi(\mathbf{z}_{n-1}, \mathbf{z}_n) z_{n-1,j} z_{nk}$$

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

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 $\xi(z_{n-1,j}, z_{nk})$ 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. θ , i.e. A , π , and ϕ_k , 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. θ , i.e. A , π , and ϕ_k , 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. θ , i.e. A , π , and ϕ_k , 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. θ , i.e. A , π , and ϕ_k , 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. θ , i.e. A , π , and ϕ_k , 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 γ and ζ

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

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

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

Computing the new parameters

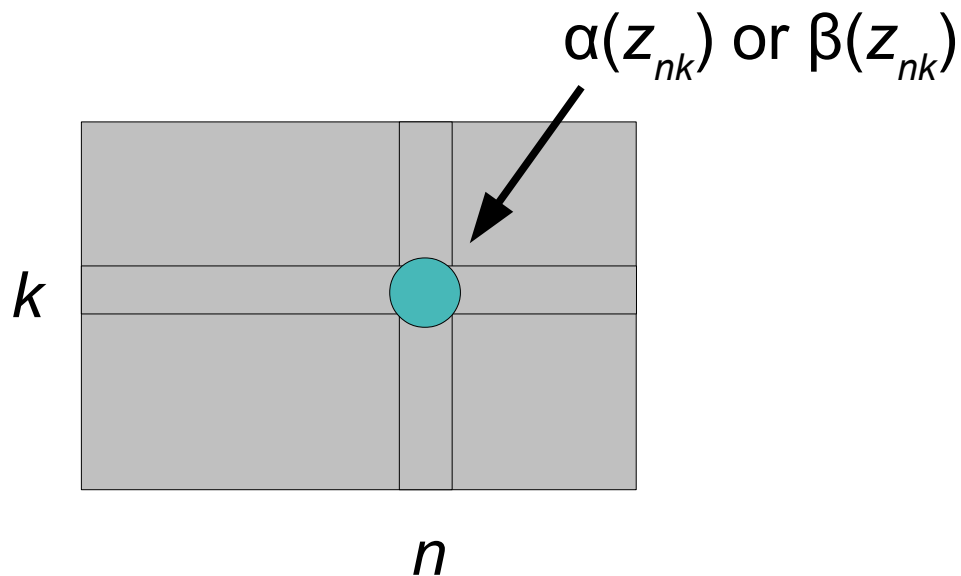
$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} = \frac{\alpha(z_{1k})\beta(z_{1k})/p(\mathbf{X})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})/p(\mathbf{X})} = \frac{\alpha(z_{1k})\beta(z_{1k})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})}$$

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

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

$$\gamma(\mathbf{z}_n) = \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

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



Computing the new parameters

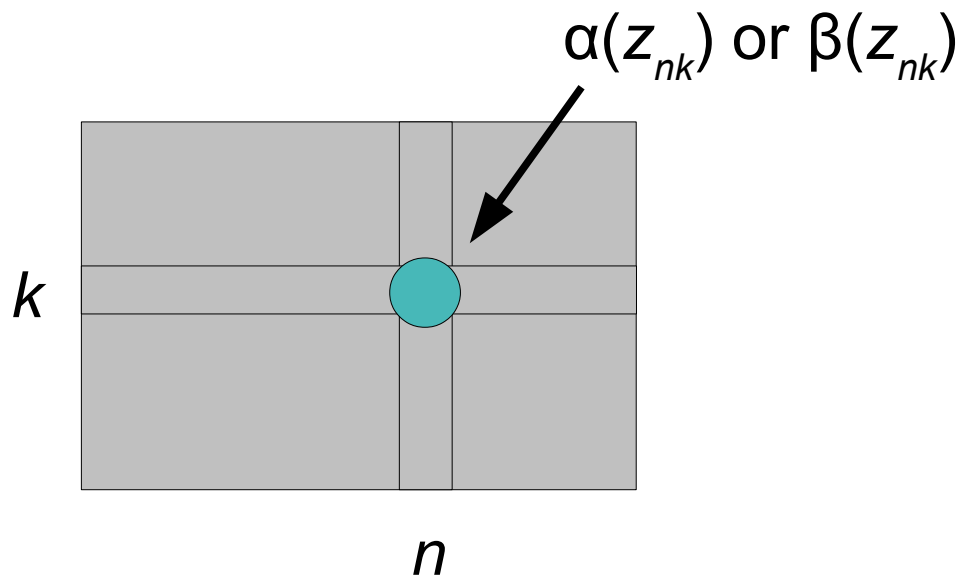
$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} = \frac{\alpha(z_{1k})\beta(z_{1k})/p(\mathbf{X})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})/p(\mathbf{X})} = \frac{\alpha(z_{1k})\beta(z_{1k})}{\sum_{j=1}^K \alpha(z_{1j})\beta(z_{1j})}$$

The old parameters

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

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

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

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

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

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

Using the scaled val

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

$$p(\mathbf{X}) = \prod_{n=1}^N c_n$$

$$\alpha(\mathbf{z}_n) = \left(\prod_{m=1}^n c_m \right) \hat{\alpha}(\mathbf{z}_n)$$

$$\beta(\mathbf{z}_n) = \left(\prod_{m=n+1}^N c_m \right) \hat{\beta}(\mathbf{z}_n)$$

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

Error in book

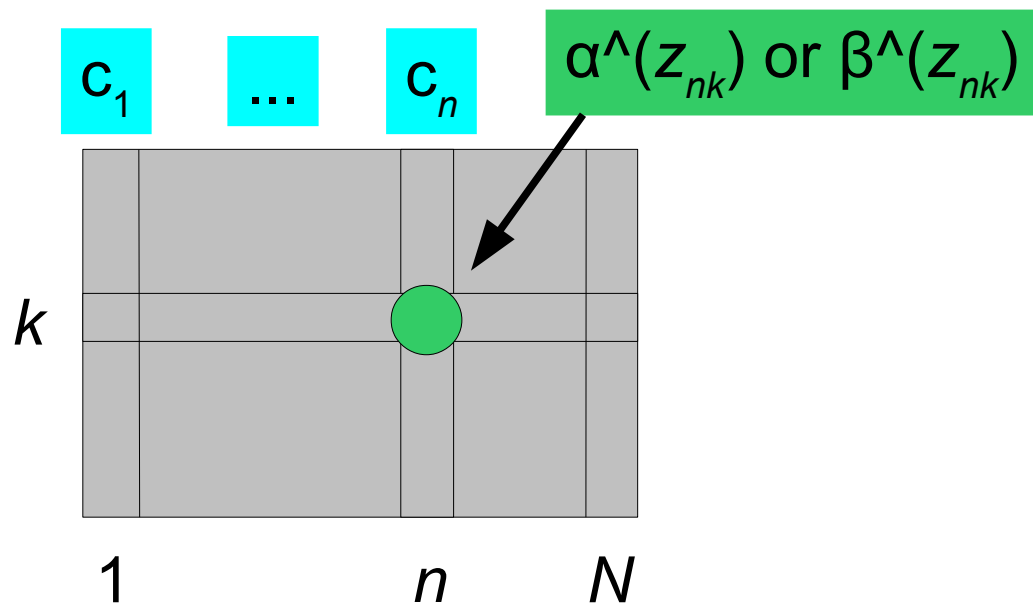
Can be computed using the modified forward- and backward-algorithm

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} ...

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