

Forward Probability , Backward Probability, Baum Welch Algorithm

HMM

HMM parameters are given by: $\Theta = (A, B, \Pi)$

HMM Definition:

- Set of states : S where $|S|=N$
- Start state : S_0 /* $P(S_0)=1$ */
- Output Alphabet : O where $|O|=K$
- State Transition Probabilities :
 $A = \{a_{ij}\}$ (prob. of going from state S_i to state S_j)
- Emission Probabilities : $B = \{b_j(O_k)\}$ (prob. Of outputting symbol O_k from state j)
- Initial State Probabilities : $\Pi = \{\Pi_i\} = \{P(X_1=i)\}$

Baum Welch Algorithm

- When corpus is available but not labeled, will use EM algorithm to estimate parameters of HMM.
- EM algorithm is used to find the maximum likelihood estimate of the parameters of HMM.

EM Algorithm

1. Choose some initial values for $\Theta = (A, B, \Pi)$.
2. Repeat the following step until convergence.
3. Determine probable (state) paths:
....., $X_{t-1} = i, X_t = j, \dots$
4. Count the expected number of transitions, $\{a_{ij}\}$, as well as the expected number of times various emissions, $b_j(O_k)$, are made.
5. Re-estimate $\Theta = (A, B, \Pi)$ using $\{a_{ij}\}$ and $b_j(O_k)$.

A forward-backward algo is used for finding probable state paths.

Forward-Backward Algorithm

Forward Probability

- $\alpha_k(i) = F(k, i) = \text{Probability of being in state } S_i \text{ after having seen } O_0 O_1 O_2 \dots O_k$
i.e., The forward variable $\alpha_k(i)$ or $F(k, i)$ is defined as the joint probability of the partial observation sequence $O_0 O_1 O_2 \dots O_k$ and that the hidden state at time k is S_i :
- $\alpha_k(i) = F(k, i) = P(O_0 O_1 O_2 \dots O_k, S_i)$
- With m as length of observation sequence

$$\begin{aligned} P(\text{Observed Sequence}) &= P(O_0 O_1 O_2 \dots O_m) \\ &= \sum_{i=0}^N P(O_0 O_1 O_2 \dots O_m, S_i) \end{aligned}$$

Forward Probability

- $P(\text{Observed Sequence}) = \sum_{i=0}^N F(m, i)$
- Method to re-compute forward recursively,

$$F(k, q) = \sum_{p=0}^N F(k-1, p) P(S_p \xrightarrow{O_k} S_q)$$

Backward Probability

- We define backward probability as:

$B(k, i)$ = Probability of seeing $O_k O_{k+1} O_{k+2} \dots O_m$ given that the state was S_i .

$$B(k, i) = P(O_k O_{k+1} O_{k+2} \dots O_m / S_i)$$

Where m is length of observed sequence.

This is better expressed by working through the mathematical calculation

Backward Probability

$$\begin{array}{ccccccc}
 O_0 & O_1 & \dots & O_{k-1} & O_k & O_{k+1} & \dots O_{m-1} & O_m \\
 S_0 \rightarrow & & & S_p & S_q & & S_m & S_{\text{final}}
 \end{array}$$

- So, the backward probability $B(k, p)$ is the probability of seeing $O_k O_{k+1} O_{k+2} \dots O_m$ given that the state was S_p .

$$B(k, p) = P(O_k O_{k+1} O_{k+2} \dots O_m / S_p)$$

$$B(k, p) = P(O_{k+1} O_{k+2} \dots O_m, O_k / S_p)$$

perform marginalization

$$B(k, p) = \sum_{q=0}^N P(O_{k+1} O_{k+2} \dots O_m, O_k, S_q / S_p)$$

- So, backward probability is :

$$\begin{aligned}
 B(k, p) &= \sum_{q=0}^N P(O_k, S_q / S_p) \cdot P(O_{k+1} O_{k+2} \dots O_m / O_k, S_q, S_p) \\
 &= \sum_{q=0}^N P(O_{k+1} O_{k+2} \dots O_m / S_q) \cdot P(O_k, S_q / S_p) \\
 &= \sum_{q=0}^N B(k+1, q) \cdot P(S_p \xrightarrow{O_k} S_q)
 \end{aligned}$$

where ,

K+1 goes on increasing towards end of observation sequence.

Boundary Condition for $B(k, p)$

- Boundary condition for $B(k, p)$ is obtained from the last symbol $B(m, \text{final})$ where S_{final} is one of the states of HMM.

Backward Probability

- $B(k, i) =$ Probability of seeing $O_k O_{k+1} O_{k+2} \dots O_m$ given that the state was S_i .
- $B(k, i) = P(O_k O_{k+1} O_{k+2} \dots O_m / S_i)$
- With m as length of observation sequence

$$\begin{aligned} P(\text{Observed Sequence}) &= P(O_0 O_1 O_2 \dots O_m) \\ &= P(O_0 O_1 O_2 \dots O_m / S_0) \\ &= B(0, 0) \end{aligned}$$

$$B(k, p) = \sum_{q=0}^N B(k+1, q) \cdot P(S_p \xrightarrow{O_k} S_q)$$

Baum Welch algorithm

- ✓ Training Hidden Markov Model (not structure learning, i.e., the structure of the HMM is pre-given). This involves:
- ✓ Learning probability values ONLY

The Learning Problem

Learning: Given an observation sequence O and the set of possible states in the HMM, learn the HMM parameters A and B .

- **Baum-Welch = Forward-Backward Algorithm**
(Baum 1972)
- Is a special case of the EM or Expectation-Maximization algorithm (Dempster, Laird, Rubin)
- The algorithm will let us train the transition probabilities $A = \{a_{ij}\}$ and the emission probabilities $B = \{b_i(o_t)\}$ of the HMM

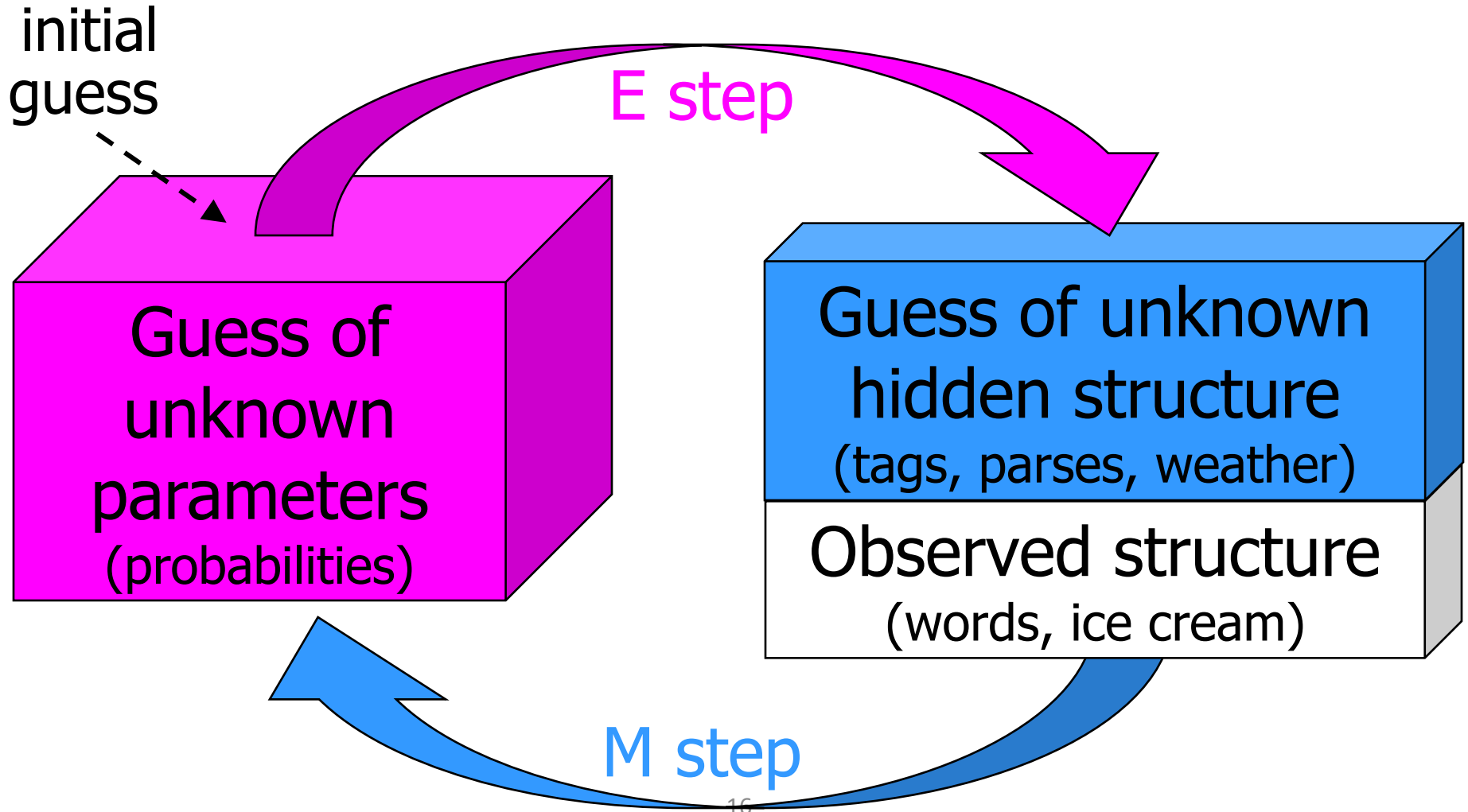
General Idea

- Start by devising a model
 - Any model that predicts the corpus observations via some hidden structure (tags, ...)
- Initially **guess** the parameters of the model!
 - Educated guess is best, but random can work

- **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure
- **Maximization step:** Use that hidden structure (and observations) to re-estimate parameters

Repeat until convergence!

General Idea



For Hidden Markov Models

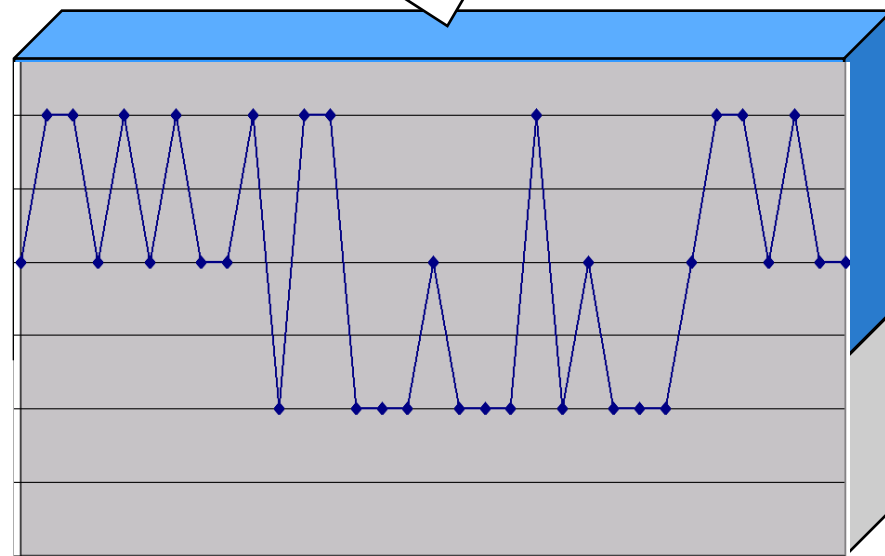
initial
guess

E step

Guess of

	$p(\dots C)$	$p(\dots H)$	$p(\dots START)$
$p(1 \dots)$	0.7	0.1	
$p(2 \dots)$	0.2	0.2	
$p(3 \dots)$	0.1	0.7	
$p(C \dots)$	0.8	0.1	0.5
$p(H \dots)$	0.1	0.8	0.5
$p(STOP \dots)$	0.1	0.1	0

(probabilities)



M step

For Hidden Markov Models

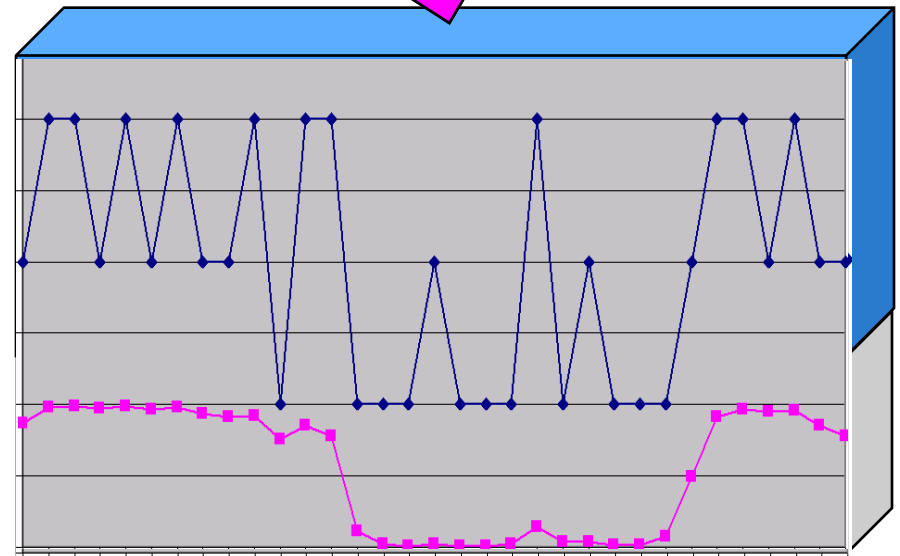
initial
guess

E step

Guess of

	$p(\dots C)$	$p(\dots H)$	$p(\dots START)$
$p(1 \dots)$	0.677	0.058	
$p(2 \dots)$	0.219	0.425	
$p(3 \dots)$	0.105	0.517	
$p(C \dots)$	0.876	0.093	0.129
$p(H \dots)$	0.109	0.865	0.871
$p(STOP \dots)$	0.015	0.042	0

(probabilities)



M step

For Hidden Markov Models

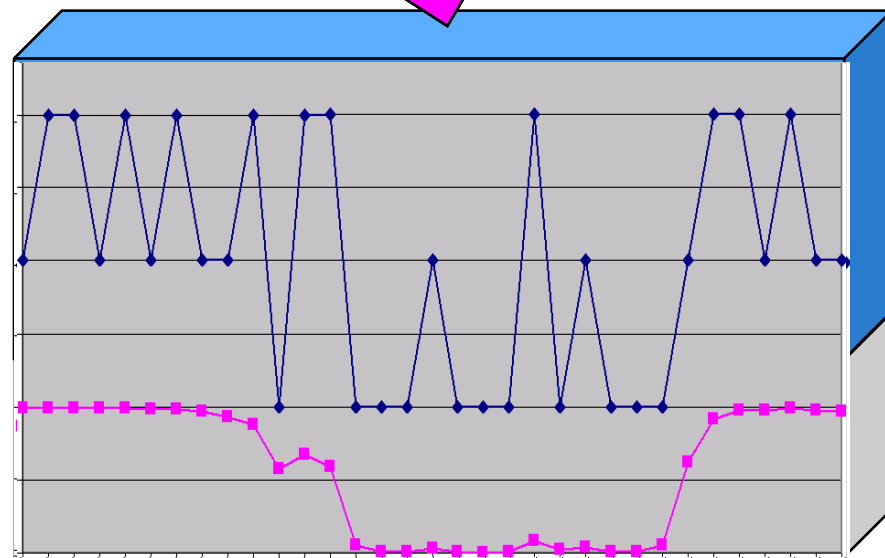
initial
guess

E step

Guess of

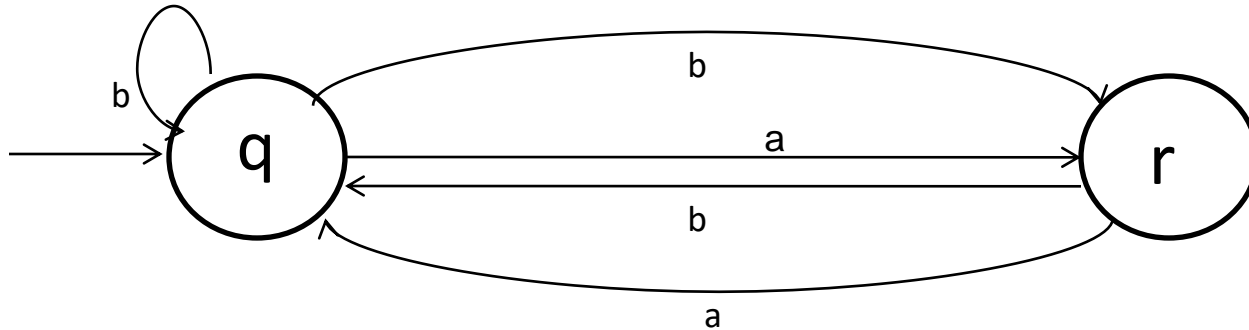
	$p(\dots C)$	$p(\dots H)$	$p(\dots START)$
$p(1 \dots)$	0.697	0.04	
$p(2 \dots)$	0.171	0.464	
$p(3 \dots)$	0.132	0.496	
$p(C \dots)$	0.904	0.077	0.012
$p(H \dots)$	0.094	0.87	0.988
$p(STOP \dots)$	0.002	0.053	0

(probabilities)



M step

Start of baum-welch algorithm



Training sequence/Observation sequence = abb aaa bbb aaa

Sequence of states with respect to input symbols

o/p seq \rightarrow $q \xrightarrow{a} r \xrightarrow{b} q \xrightarrow{b} q \xrightarrow{a} r \xrightarrow{a} q \xrightarrow{a} r \xrightarrow{b} q \xrightarrow{b} q \xrightarrow{b} q \xrightarrow{a} r \xrightarrow{a} q \xrightarrow{a} r$
State seq

Calculating probabilities from table

$$P(q \xrightarrow{a} r) = 5/8$$

$$P(q \xrightarrow{b} q) = 3/8$$

$$P(s^i \xrightarrow{w_k} s^j) = \frac{c(s^i \xrightarrow{w_k} s^j)}{\sum_{l=1}^T \sum_{m=1}^A c(s^i \xrightarrow{w_m} s^l)}$$

$T = \#states$

$A = \#alphabet\ symbols$

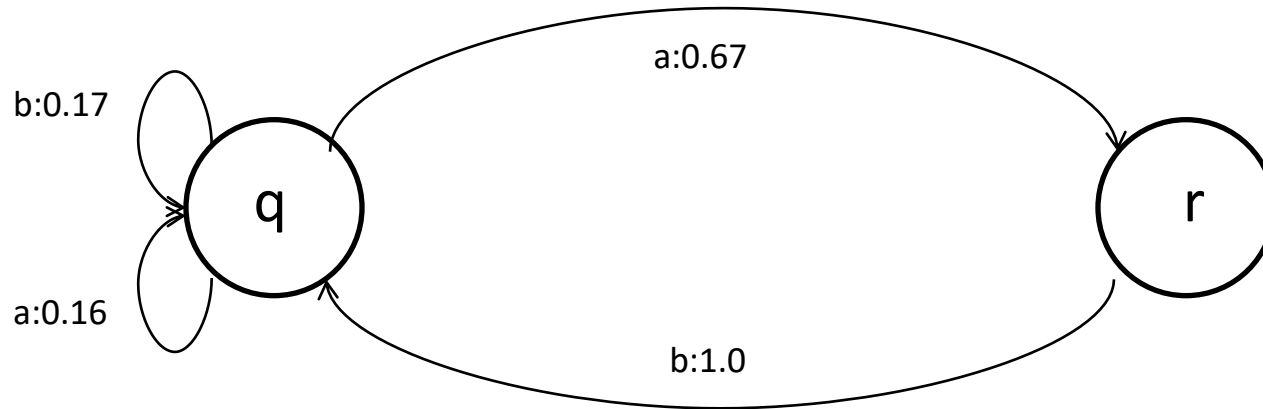
Now if we have a non-deterministic transitions then multiple state seq possible for the given o/p seq (ref. to previous slide's feature). Our aim is to find expected count through this.

Table of counts

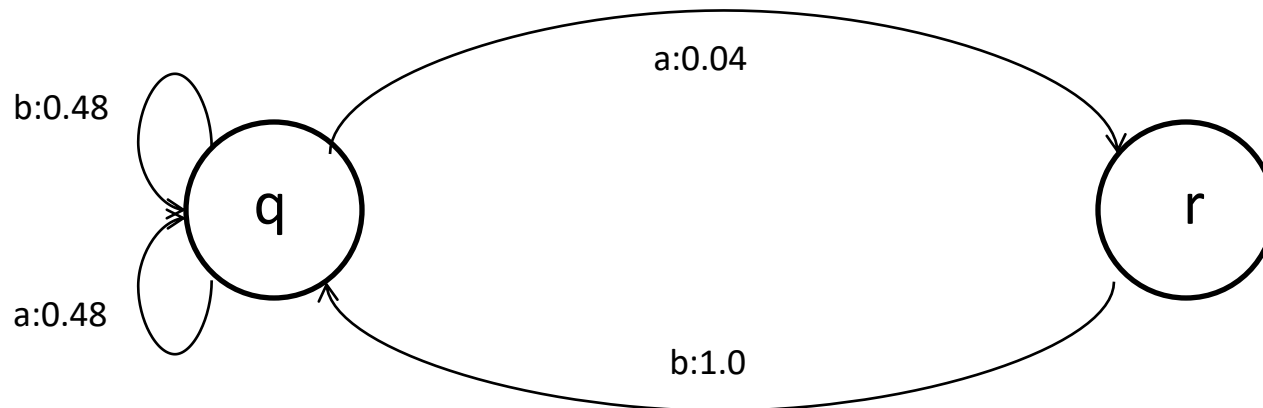
Src	Dest	O/P	Count
q	r	a	5
q	q	b	3
r	q	a	3
r	q	b	2

Learning probabilities

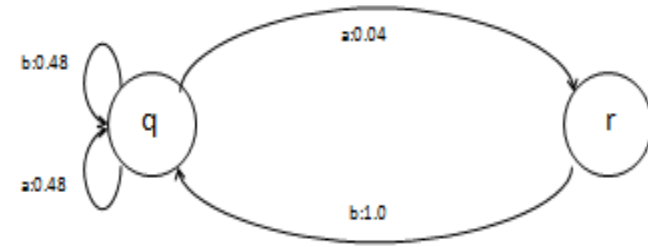
Machine Given : Learn from it



Initial guess



One run of Baum-Welch algorithm: *string ababb*



$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077				
q	r	q	q	q	q					
q	q	q	r	q	q					
q	q	q	q	q	q					
Expected Count \rightarrow										
New Probabilities (P) \rightarrow										

State sequences

* ϵ is considered as starting and ending symbol of the input sequence string

$P(\text{state seq: qrqrqq}) =$

$$\begin{aligned}
 P(\text{qrqrqq}) &= P(S_0 \xrightarrow{\epsilon} q) * P(q \xrightarrow{a} r) * P(r \xrightarrow{b} q) * P(q \xrightarrow{a} r) * P(r \xrightarrow{b} q) * P(q \xrightarrow{b} q) \\
 &= 1 * 0.04 * 1.0 * 0.04 * 1.0 * 0.48 \\
 &= 0.00077
 \end{aligned}$$

One run of Baum-Welch algorithm: *string ababb*

$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077				
q	r	q	q	q	q	0.00442				
q	q	q	r	q	q	0.00442				
q	q	q	q	q	q	0.02548				
Expected Count \rightarrow										
New Probabilities (P) \rightarrow										

State sequences

* ϵ is considered as starting and ending symbol of the input sequence string

$$\begin{aligned}
 P(qrqrqq) &= P(S_0 \xrightarrow{\epsilon} q) * P(q \xrightarrow{a} r) * P(r \xrightarrow{b} q) * P(q \xrightarrow{a} r) * P(r \xrightarrow{b} q) * P(q \xrightarrow{b} q) \\
 &= 1 * 0.04 * 1.0 * 0.04 * 1.0 * 0.48 \\
 &= 0.00077
 \end{aligned}$$

One run of Baum-Welch algorithm: *string ababb*

$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077				
q	^r ↑ q	q	q	q	q	0.00442				
q	q	q	r	q	q	0.00442				
q	q	q	q	q	q	0.02548				
Expected Count →						0.035				
New Probabilities (P) →										

State sequences

* ϵ is considered as starting and ending symbol of the input sequence string

$$\begin{aligned}
 \text{Expected Count (path)} &= \text{Summation(path Column)} \\
 &= (0.00077 + 0.00442 + 0.00442 + 0.2548) \\
 &= 0.03509 \approx \mathbf{0.035}
 \end{aligned}$$

One run of Baum-Welch algorithm:

string ababb

$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077	0.00154			
q	r	q	q	q	q	0.00442				
q	q	q	r	q	q	0.00442				
q	q	q	q	q	q	0.02548				
Expected Count \rightarrow						0.035				
New Probabilities (P) \rightarrow										

Count for state transitions, i.e., $q \xrightarrow{a} r$, $r \xrightarrow{b} q$, and so on.

$$\text{Count}(q \xrightarrow{a} r) = \text{path} * \text{number of transitions} = 0.00077 * 2 = 0.00154$$

One run of Baum-Welch algorithm:

string ababb

$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077	0.00154	0.00154	0	0.00077
q	r	q	q	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	r	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	q	q	q	0.02548	0.0	0.000	0.05096	0.07644
Expected Count \rightarrow						0.035	0.01	0.01	0.06	0.095
New Probabilities (P) \rightarrow										

One run of Baum-Welch algorithm: *string ababb*

$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077	0.00154	0.00154	0	0.00077
q	r	q	q	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	r	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	q	q	q	0.02548	0.0	0.000	0.05096	0.07644
Expected Count \rightarrow						0.035	0.01	0.01	0.06	0.095
New Probabilities (P) \rightarrow							0.06			

Estimating new probabilities for state sequences, i.e., $q \xrightarrow{a} r$ $r \xrightarrow{b} q$ and so on.

$$\begin{aligned}
 \text{Normalized new probability for } q \xrightarrow{a} r &= \frac{\text{expected count of } q \xrightarrow{a} r}{\sum \text{expected count of } q \xrightarrow{*} *} \\
 &= \frac{0.01}{0.01 + 0.06 + 0.095} \\
 &= 0.06
 \end{aligned}$$

One run of Baum-Welch algorithm:

string ababb

$\epsilon \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow \epsilon$	P(path)	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077	0.00154	0.00154	0	0.00077
q	r	q	q	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	r	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	q	q	q	0.02548	0.0	0.000	0.05096	0.07644
Expected Count \rightarrow						0.035	0.01	0.01	0.06	0.095
New Probabilities (P) \rightarrow							0.06	1.0	0.36	0.581

State sequences

Note:

- Now new path values can be recomputed.
- So new expected counts which will lead to new probability values.
- This way through multiple iterations the probability values will converge.

What is an epoch?

- It is one iteration over all observation patterns.
- So, after each iteration we get updated probabilities. This is done until converges.

The Baum-Welch Algorithm

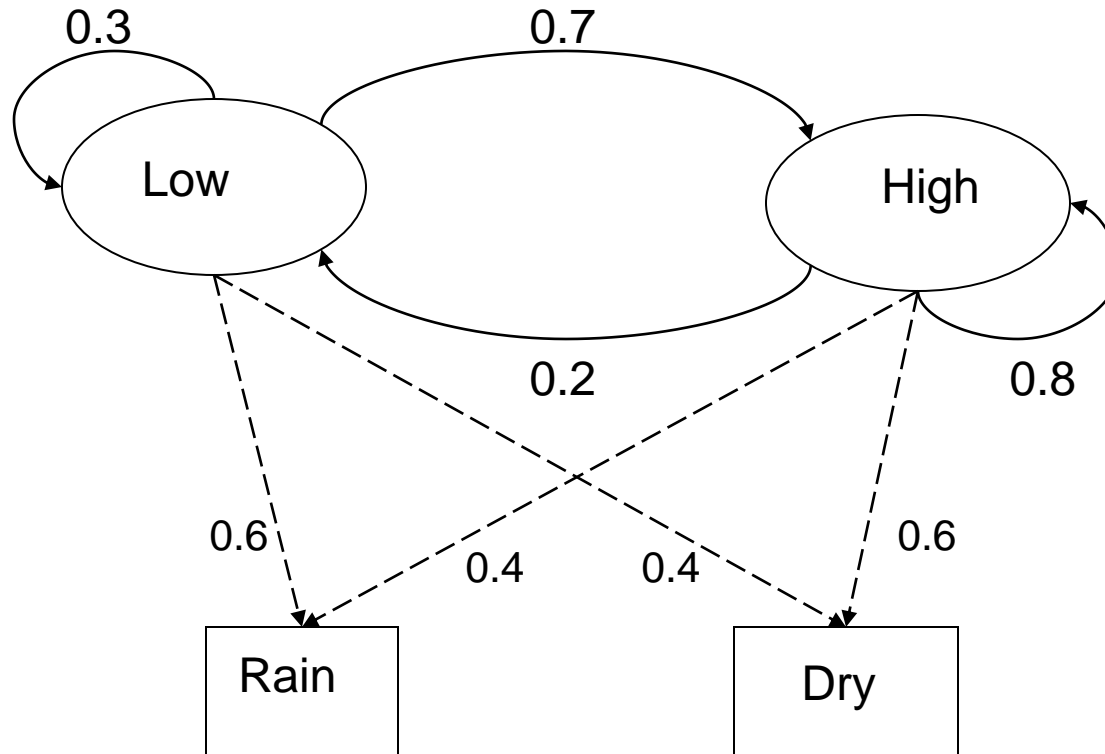
- initialize the parameters of the HMM
- iterate until convergence
 - initialize $n_{k,c}$, $n_{k \rightarrow l}$ with pseudocounts
 - **E-step**: for each training set sequence $j = 1 \dots n$
 - calculate $f_k(i)$ values for sequence j
 - calculate $b_k(i)$ values for sequence j
 - add the contribution of sequence j to $n_{k,c}$, $n_{k \rightarrow l}$
 - **M-step**: update the HMM parameters using $n_{k,c}$, $n_{k \rightarrow l}$

where,

$n_{k,c}$ be the expected number of emissions of c from state k for the training set, and

$n_{k \rightarrow l}$ be the expected number of transitions from state k to state l for the training set

Example of Hidden Markov Model



Example of Hidden Markov Model

1. Two states : 'Low' and 'High' atmospheric pressure.
2. Two observations : 'Rain' and 'Dry'.
3. Transition probabilities: $P(\text{'Low'}|\text{'Low'})=0.3$, $P(\text{'High'}|\text{'Low'})=0.7$, $P(\text{'Low'}|\text{'High'})=0.2$, $P(\text{'High'}|\text{'High'})=0.8$
4. Observation probabilities :
 $P(\text{'Rain'}|\text{'Low'})=0.6$, $P(\text{'Dry'}|\text{'Low'})=0.4$,
 $P(\text{'Rain'}|\text{'High'})=0.4$, $P(\text{'Dry'}|\text{'High'})=0.3$.
5. Initial probabilities: say $P(\text{'Low'})=0.4$,
 $P(\text{'High'})=0.6$.

Calculation of observation sequence probability

1. Suppose we want to calculate a probability of a sequence of observations in our example, {'Dry','Rain'}.
2. Consider all possible hidden state sequences:
3. $P(\text{'Dry','Rain'}) = P(\text{'Dry','Rain'} , \text{'Low','Low'}) + P(\text{'Dry','Rain'} , \text{'Low','High'}) + P(\text{'Dry','Rain'} , \text{'High','Low'}) + P(\text{'Dry','Rain'} , \text{'High','High'})$

where first term is :

$$\begin{aligned} &P(\text{'Dry','Rain'} , \text{'Low','Low'}) = P(\text{'Dry','Rain'} | \text{'Low','Low'}) P(\text{'Low','Low'}) = \\ &P(\text{'Dry'}|\text{'Low'})P(\text{'Rain'}|\text{'Low'}) P(\text{'Low'})P(\text{'Low'}|\text{'Low'}) \\ &= 0.4*0.4*0.6*0.4*0.3 \end{aligned}$$