



Royal Institute of Technology

MACHINE LEARNING 2 – EM, HMMS (CH 17)

Lecture 4

RELATIONS BETWEEN LOG-LIKELIHOODS AND Q-TERMS

Theorem: for $\theta' = \operatorname{argmax}_{\theta'} Q(\theta', \theta)$

$$\log p(\mathcal{D}|\theta') \geq Q(\theta', \theta) - R(\theta, \theta) \geq Q(\theta, \theta) - R(\theta, \theta) = \log p(\mathcal{D}|\theta)$$

So by maximizing Q-term (through ESS) we monotonically increase the likelihood.

The Q-term may not increase in every step!

Expected complete: of all data

$\log p(\mathcal{D}|\theta')$

$$\begin{aligned} &= \sum_n \log \sum_{z_n} p(\mathbf{x}_n, z_n|\theta') \\ &= \sum_n \log \sum_{z_n} p(z_n|\mathbf{x}_n, \theta) \frac{p(\mathbf{x}_n, z_n|\theta')}{p(z_n|\mathbf{x}_n, \theta)} \\ &\geq \sum_n \sum_{z_n} p(z_n|\mathbf{x}_n, \theta) \log \frac{p(\mathbf{x}_n, z_n|\theta')}{p(z_n|\mathbf{x}_n, \theta)} \\ &= \sum_n \sum_{z_n} p(z_n|\mathbf{x}_n, \theta) \log p(\mathbf{x}_n, z_n|\theta') - \sum_{z_n} p(z_n|\mathbf{x}_n, \theta) \log p(z_n|\mathbf{x}_n, \theta) \\ &= \underbrace{\sum_n Q_n(\theta'; \theta)}_{Q(\theta'; \theta)} - \underbrace{\sum_n R_n(\theta; \theta)}_{R_n(\theta; \theta)} \end{aligned}$$

GMM EXPECTED LOG LIKELIHOOD

$$\begin{aligned} E_{p(z_n|\mathbf{x}_n, \theta)} [l(\theta'; \mathcal{D})] &= E_{p(z_n|\mathbf{x}_n, \theta)} \left[\log \prod_n p(\mathbf{x}_n, z_n|\theta') \right] \\ &= \sum_n E \left[\log \prod_c (\pi'_c p(\mathbf{x}_n|z_n=c, \theta'))^{I(z_n=c)} \right] \\ &= \sum_n \sum_c E [I(z_n=c) \log(\pi'_c p(\mathbf{x}_n|\theta'_c))] \\ &= \sum_n \sum_c p(z_n=c|\mathbf{x}_n, \theta) \log(\pi'_c p(\mathbf{x}_n|\theta'_c)) \\ &= \sum_n \sum_c r_{nc} \log \pi'_c + \sum_n \sum_c r_{nc} \log p(\mathbf{x}_n|\theta'_c) \end{aligned}$$

where for a one dim Gaussian $\theta_c = (\mu_c, \sigma_c^2)$

DGM WITH CATEGORICAL CPDS - LIKELIHOOD

$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

- Let $\theta'_{vkl} = p(X_v = k | \mathbf{X}_{\text{pa}(v)} = \mathbf{l}, \theta')$, then

$$p(X_v = k | \mathbf{X}_{\text{pa}(v)} = \mathbf{l}, \theta') = \prod_{k,l} (\theta'_{vkl})^{I(X_v=k, \mathbf{X}_{\text{pa}(v)}=\mathbf{l})}$$

- The loglikelihood $\log p(\mathcal{D} | \theta')$ equals

$$\sum_n \sum_{v,k,l} I(x_{nv} = k, \mathbf{x}_{n,\text{pa}(v)} = \mathbf{l}) \log \theta'_{v,k,l} = \sum_{v,k,l} N_{v,k,l} \log \theta'_{v,k,l}$$

where $N_{v,k,l} = \sum_n I(x_{nv} = k, \mathbf{x}_{n,\text{pa}(v)} = \mathbf{l})$

EXPECTED COMPLETE LOG-LIKELIHOOD

- The expected complete log-likelihood is

$$\sum_n \sum_{v,k,l} E_{p(\mathbf{X} | \mathbf{x}_n, \theta)} [I(X_v = k, \mathbf{X}_{\text{pa}(v)} = \mathbf{l})] \log \theta'_{v,k,l} = \sum_{v,k,l} \bar{N}_{v,k,l} \log \theta'_{v,k,l}$$

where

$$\begin{aligned} \bar{N}_{v,k,l} &= \sum_n E_{p(\mathbf{X} | \mathbf{x}_n, \theta)} [I(X_v = k, \mathbf{X}_{\text{pa}(v)} = \mathbf{l})] && \text{Can be computed using GM inference} \\ &= \sum_n p(X_v = k, \mathbf{X}_{\text{pa}(v)} = \mathbf{l} | \mathbf{x}_n, \theta) \end{aligned}$$

- We can independently maximize each

$$\sum_k \bar{N}_{v,k,l} \log \theta'_{v,k,l}$$

- Done by setting

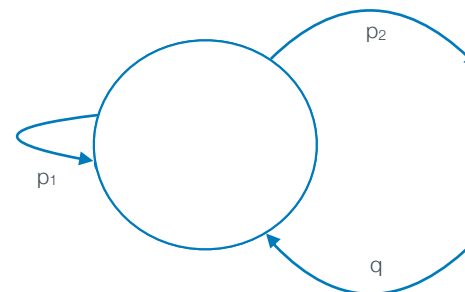
$$\theta'_{v,k,l} = \frac{\bar{N}_{v,k,l}}{\sum_k \bar{N}_{v,k,l}}$$

MAXIMIZED
INDEPENDENTLY

MARKOV CHAINS (DISCRETE)

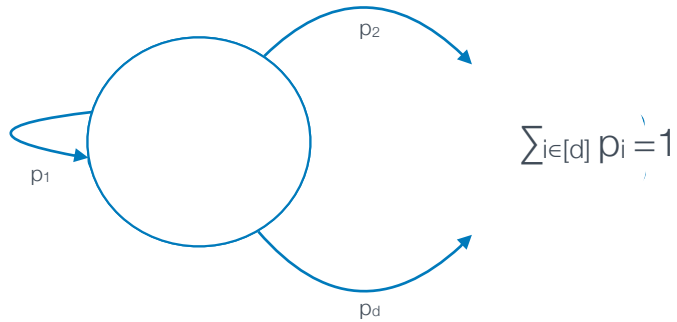
★ Directed graph with transition probabilities

★ We observe the sequence of visited vertices



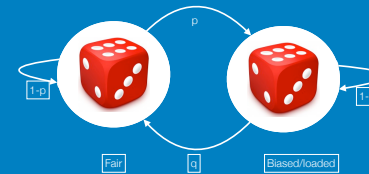
MARKOV CHAINS (DISCRETE)

Probabilities on outgoing edges sum to one



WHAT AN HMM DOES

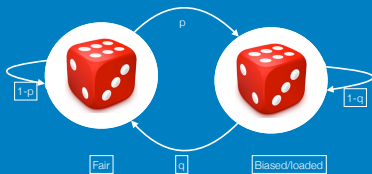
Rolls: 6641532161621152346532143566342616552
Die: LLLLLLLLLLLLLLFFFFFLLLLLLLLLLLLLLLLLFFF



- ★ Starts in the state z_1
- ★ When in state z_t
 - outputs $p(x_t|z_t)$
 - moves to $p(z_{t+1}|z_t)$
- ★ Stops after a fixed number of steps or when reaching a stop step

WHAT AN HMM DOES

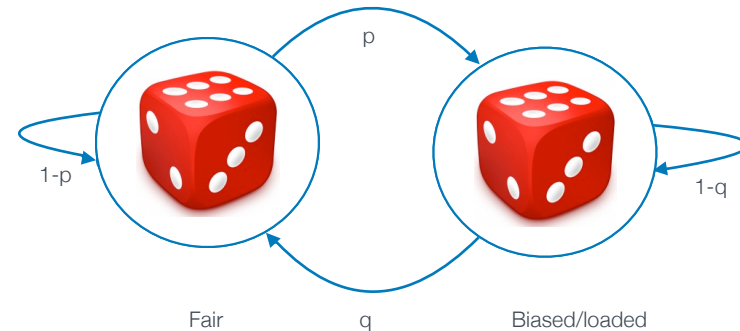
Rolls: 6641532161621152346532143566342616552
Die: LLLLLLLLLLLLLLFFFFFLLLLLLLLLLLLLLLLLFFF



- ★ Starts in the state z_1
- ★ When in state z_t
 - outputs $p(x_t|z_t)$
 - moves to $p(z_{t+1}|z_t)$
- ★ Stops after a fixed number of steps or when reaching a stop step

↑
↑
The parameters

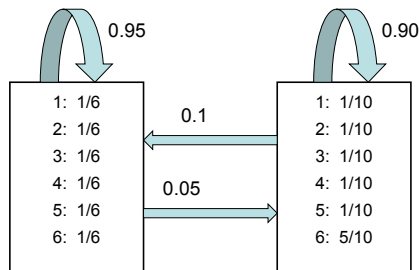
THE OCCASIONALLY DISHONEST CASINO



- ★ We observe the sequence of dice outcomes of visited vertices

Rolls: 664153216162115234653214356634261655234232315142464156663246
Die: LLLLLLLLLLLLLLFFFFFLLLLLLLLLLLLLLLLLFFF

EMISSION DISTRIBUTIONS



THE JOINT DISTRIBUTION

$$\begin{aligned}
 p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) &= p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}) p(\mathbf{z}_{1:T}) \\
 &= p(\mathbf{z}_1) \left(\prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) \right) \left(\prod_{t=1}^T p(\mathbf{x}_t | \mathbf{z}_t) \right)
 \end{aligned}$$

Categorical or Gaussian

- ★ Starts in the state z_1
- ★ When in state z_t
 - emits $p(x_t | z_t)$
 - transits to $p(z_{t+1} | z_t)$
- ★ Stops after a fixed number of steps or when reaching a stop step

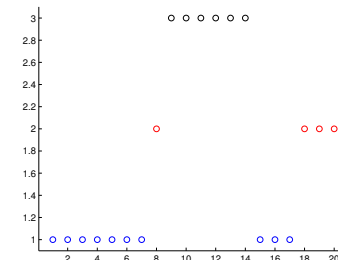
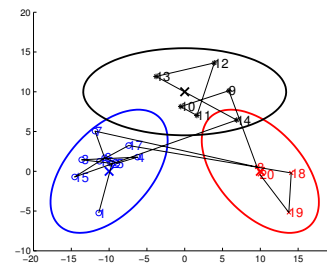
THE JOINT DISTRIBUTION

$$\begin{aligned}
 p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}) &= p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}) p(\mathbf{z}_{1:T}) \\
 &= p(\mathbf{z}_1) \left(\prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) \right) \left(\prod_{t=1}^T p(\mathbf{x}_t | \mathbf{z}_t) \right)
 \end{aligned}$$

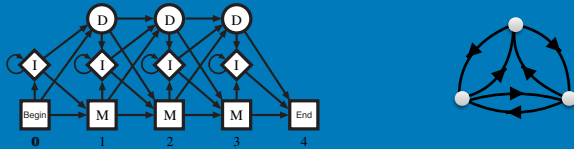
Categorical or Gaussian

- ★ Starts in the state z_1
- ★ When in state z_t
 - emits $p(x_t | z_t)$
 - transits to $p(z_{t+1} | z_t)$
- ★ Stops after a fixed number of steps or when reaching a stop step

GAUSSIAN EMISSIONS AND HIDDEN STATES

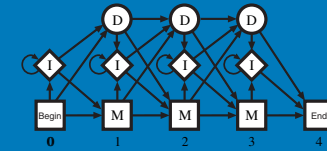


LAYERED OR NOT



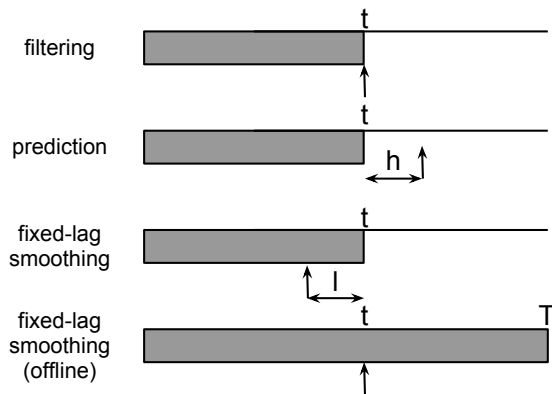
APPLICATIONS OF HMMS

	x	x	.	.	x
bat	A	G	-	-	C
rat	A	-	A	G	-
cat	A	G	-	A	A
gnat	-	-	A	A	A
goat	A	G	-	-	C
	1	2	.	.	3



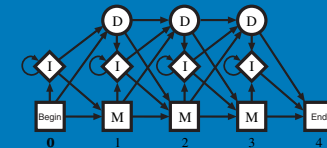
- Automatic speech recognition
- Activity recognition
- Part of speech tagging
- Gene finding
- Protein sequence alignment

TERMINOLOGY X ABOVE Z BELOW



INFERENCE TYPES

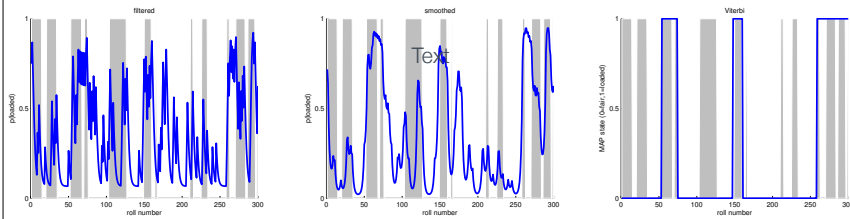
	x	x	.	.	x
bat	A	G	-	-	C
rat	A	-	A	G	-
cat	A	G	-	A	A
gnat	-	-	A	A	A
goat	A	G	-	-	C
	1	2	.	.	3



- Filtering: $p(z_t|x_{1:t})$, online
- Smoothing (MAP): $p(z_t|x_{1:T})$ offline
- Fixed lag smoothing: $p(z_t|x_{1:t+h})$
- Prediction: $p(z_{t+h}|x_{1:t})$
- Viterbi (MAP) $\text{argmax } p(z_{1:T}|x_{1:T})$
- Posterior samples: $\sim p(z_{1:T}|x_{1:T})$
- Probability of data: $p(x_{1:T})$
- Parameters: given D & struct.
- Structure and param.: given D

INFERENCE IN THE OCCASIONALLY DISHONEST CASINO

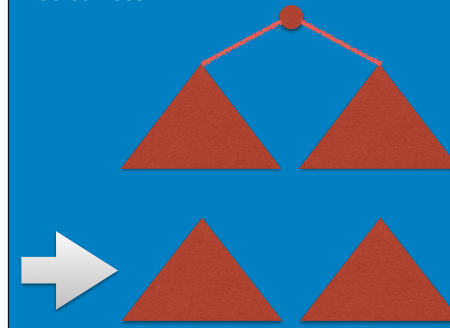
Grey regions are states corresponding to biased die



- Filtering: $p(z_t|x_{1:t})$, online
- Smoothing, MAP state: $p(z_t|x_{1:T})$ offline
- Viterbi, MAP path $\text{argmax } p(z_{1:T}|x_{1:T})$

Pairs of strings abbacd acbadd
 abbac acbadd
 abbacd acbad
 abbac acbad

Rooted trees



DP

- What is a subproblem?
- What is a subsolution?
- How do we decompose into smaller subproblems?
- How do we combine subsolutions into larger?
- How do we enumerate?
- How many and what time?

DP

Polynomial many

Polynomial time

Polynomial time

Polynomial time

Polynomial time overall

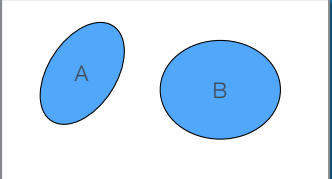
- What is a subproblem?
- What is a subsolution?
- How do we decompose into smaller subproblems?
- How do we combine subsolutions into larger?
- How do we enumerate?
- How many and what time?

SPECIAL CASE: HIDDEN MARKOV MODEL (HMM)

$$\begin{array}{cccc}
 z_1 & \rightarrow & z_2 & \rightarrow & z_3 & \cdots & z_V \\
 \downarrow & & \downarrow & & \downarrow & & \downarrow \\
 x_1 & & x_2 & & x_3 & & x_V
 \end{array}$$

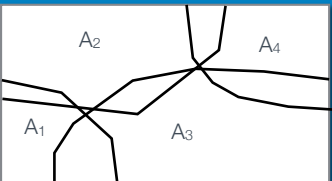
- Z_i hidden
- X_i observable
- Hidden often not observable when training, never when applying

Exclusive




EXCLUSIVE & EXHAUSTIVE

Exhaustive



Exclusive & exhaustive



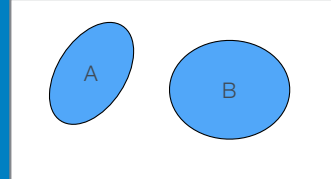
- Exclusive

$$p(A \text{ or } B) = p(A) + P(B)$$

- Exclusive & exhaustive

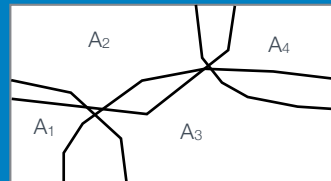
$$\sum_i p(A_i) = 1$$

Exclusive

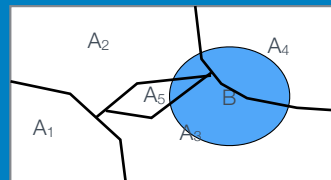


EXCLUSIVE & EXHAUSTIVE

Exhaustive



Exclusive & exhaustive



- Exclusive

$$p(A \text{ or } B) = p(A) + P(B)$$

- Exclusive & exhaustive

$$p(B) = \sum_i p(B, A_i) = \sum_i p(A_i)p(B|A_i)$$

$$p(\mathbf{x}_{1:T}) = \sum_{\mathbf{z}_{1:T}} p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T})$$

$$f_t(k) := p(\mathbf{x}_{1:t-1}, \mathbf{Z}_t = k)$$

FORWARD

- Joint is easy to express
- The sum has exponentially many terms
- The forward variable, f_t , can be computed with DP

$$p(\mathbf{x}_{1:T}) = \sum_{\mathbf{z}_{1:T}} p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T})$$

$$f_t(k) := p(\mathbf{x}_{1:t-1}, \mathbf{Z}_t = k)$$

FORWARD

- Joint is easy to express
- The sum has exponentially many terms
- The forward variable, f_t , can be computed with DP

$$\begin{aligned}
 f_t(k) &= \sum_l p(\mathbf{x}_{1:t-1}, \mathbf{Z}_{t-1} = l, \mathbf{Z}_t = k) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l) \\
 &= \sum_l \underbrace{f_{t-1}(l)}_{\text{smaller}} \underbrace{p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1} = l)}_{\text{emission}} \underbrace{p(\mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l)}_{\text{transition}}
 \end{aligned}$$

Zoom in

$$\begin{aligned}
 f_t(k) &= \sum_l p(\mathbf{x}_{1:t-1}, \mathbf{Z}_{t-1} = l, \mathbf{Z}_t = k) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l) \\
 &= \sum_l \underbrace{f_{t-1}(l)}_{\text{smaller}} \underbrace{p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1} = l)}_{\text{emission}} \underbrace{p(\mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l)}_{\text{transition}}
 \end{aligned}$$

Zoom in

K - states
T - length of observation

$$\begin{aligned}
 \begin{matrix} \nearrow \\ \nearrow \end{matrix} f_t(k) \\
 \begin{matrix} \nearrow \\ \nearrow \end{matrix} \begin{matrix} T \\ K \end{matrix} &= \sum_l p(\mathbf{x}_{1:t-1}, \mathbf{Z}_{t-1} = l, \mathbf{Z}_t = k) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l) \\
 &= \sum_l \underbrace{f_{t-1}(l)}_{\text{smaller}} \underbrace{p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1} = l)}_{\text{emission}} \underbrace{p(\mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l)}_{\text{transition}}
 \end{aligned}$$

Total time $O(TK^2)$

Zoom in

K - states
T - length of observation

$$\begin{aligned}
 \begin{matrix} \nearrow \\ \nearrow \end{matrix} f_t(k) \\
 \begin{matrix} \nearrow \\ \nearrow \end{matrix} \begin{matrix} T \\ K \end{matrix} &= \sum_l p(\mathbf{x}_{1:t-1}, \mathbf{Z}_{t-1} = l, \mathbf{Z}_t = k) \\
 \text{Constant} &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) \\
 &= \sum_l p(\mathbf{x}_{1:t-2}, \mathbf{Z}_{t-1} = l) p(\mathbf{x}_{t-1}, \mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l) \\
 &= \sum_l \underbrace{f_{t-1}(l)}_{\text{smaller}} \underbrace{p(\mathbf{x}_{t-1} | \mathbf{Z}_{t-1} = l)}_{\text{emission}} \underbrace{p(\mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l)}_{\text{transition}}
 \end{aligned}$$

Total time $O(TK^2)$

If layered, $O(TK)$

$$p(\mathbf{x}_{1:T}) = \sum_k p(\mathbf{x}_{1:T}, z_T = k)$$

$$f_t(k) = p(\mathbf{x}_{1:t-1}, \mathbf{Z}_t = k)$$

In general, (e.g. $t=T$)

$$p(\mathbf{x}_{1:t}) = \sum_k f_t(k) p(\mathbf{x}_t | \mathbf{Z}_t = k)$$

OBSERVATION PROBABILITY

- The final probability is easily obtained

FILTERING

$$p(\mathbf{Z}_t = k | \mathbf{x}_{1:t})$$

- Filtering: $p(\mathbf{z}_t | \mathbf{x}_{1:t})$, online

FILTERING

$$\begin{aligned} p(\mathbf{Z}_t = k | \mathbf{x}_{1:t}) &= \frac{p(\mathbf{x}_{1:t}, \mathbf{Z}_t = k)}{p(\mathbf{x}_{1:t})} \\ &= \frac{p(\mathbf{x}_{1:t-1}, \mathbf{Z}_t = k) p(\mathbf{x}_t | \mathbf{Z}_t = k)}{p(\mathbf{x}_{1:t})} \\ &= \frac{f_t(k) p(\mathbf{x}_t | \mathbf{Z}_t = k)}{p(\mathbf{x}_{1:t})} \end{aligned}$$

← emission
← data probability

- Filtering: $p(\mathbf{z}_t | \mathbf{x}_{1:t})$, online

BACKWARDS

$$b_t(k) := p(\mathbf{x}_{t+1:T} | \mathbf{Z}_t = k)$$

- DP also for the backward variable b_t

BACKWARDS

$$b_t(k) := p(\mathbf{x}_{t+1:T} | \mathbf{Z}_t = k)$$

- DP also for the backward variable b_t

$$\begin{aligned} b_t(k) &= \sum_l p(\mathbf{x}_{t+1:T}, \mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) \\ &= \sum_l p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) p(\mathbf{x}_{t+1:T} | \mathbf{Z}_{t+1} = l) \\ &= \sum_l p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) p(\mathbf{x}_{t+2:T} | \mathbf{Z}_{t+1} = l) p(\mathbf{x}_{t+1} | \mathbf{Z}_{t+1} = l) \\ &= \sum_l \underbrace{p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k)}_{\text{transition}} \underbrace{b_{t+1}(l)}_{\text{"smaller"}} \underbrace{p(\mathbf{x}_{t+1} | \mathbf{Z}_{t+1} = l)}_{\text{emission}} \end{aligned}$$

BACKWARDS

$$b_t(k) := p(\mathbf{x}_{t+1:T} | \mathbf{Z}_t = k)$$

- DP also for the backward variable b_t

$$\begin{aligned}
 b_t(k) &= \sum_l p(\mathbf{x}_{t+1:T}, \mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) && \text{Running time analysis as before} \\
 &= \sum_l p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) p(\mathbf{x}_{t+1:T} | \mathbf{Z}_{t+1} = l) \\
 &= \sum_l p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) p(\mathbf{x}_{t+2:T} | \mathbf{Z}_{t+1} = l) p(\mathbf{x}_{t+1} | \mathbf{Z}_{t+1} = l) \\
 &= \sum_l \underbrace{p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k)}_{\text{transition}} \underbrace{b_{t+1}(l)}_{\text{"smaller"}} \underbrace{p(\mathbf{x}_{t+1} | \mathbf{Z}_{t+1} = l)}_{\text{emission}}
 \end{aligned}$$

SMOOTHING

$$p(\mathbf{Z}_t = k | \mathbf{x}_{1:T})$$

- Smoothing, MAP state: $p(z_t | \mathbf{x}_{1:T})$, offline

SMOOTHING

$$\begin{aligned}
 p(\mathbf{Z}_t = k | \mathbf{x}_{1:T}) &= \frac{p(\mathbf{x}_{1:T}, \mathbf{Z}_t = k)}{p(\mathbf{x}_{1:T})} \\
 &\propto p(\mathbf{x}_{1:t-1}, \mathbf{Z}_t = k) p(\mathbf{x}_{t:T} | \mathbf{Z}_t = k) \\
 &= f_t(k) \underbrace{p(\mathbf{x}_t | \mathbf{Z}_t = k)}_{\text{emission}} b_t(k)
 \end{aligned}$$

- Smoothing, MAP state: $p(z_t | \mathbf{x}_{1:T})$, offline

TWO SLICED SMOOTH MARGINALS - MARGINAL OVER PAIRS OF STATES

$$p(\mathbf{Z}_t = k, \mathbf{Z}_{t+1} = l | \mathbf{x}_{1:T})$$

- Can be computed from forward and backward similarly

We want $\operatorname{argmax}_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})$

Not!

$(\operatorname{argmax}_{z_1} p(z_1 | \mathbf{x}_{1+1:T}), \dots, \operatorname{argmax}_{z_T} p(z_T | \mathbf{x}_{1+1:T}))$

Viterbi variable

$$v_t(k) := \max_{z_{1:t-1}} p(\mathbf{z}_{1:t-1}, \mathbf{Z}_t = k, \mathbf{x}_{1:t})$$

Gives what we want

VITERBI

- MAP path
- Viterbi learning: used, as approximation, to speed up parameter learning
- Again DP now with Viterbi variable
- For the path, use back pointers

We want $\operatorname{argmax}_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})$

Not!

$(\operatorname{argmax}_{z_1} p(z_1 | \mathbf{x}_{1+1:T}), \dots, \operatorname{argmax}_{z_T} p(z_T | \mathbf{x}_{1+1:T}))$

Viterbi variable

$$v_t(k) := \max_{z_{1:t-1}} p(\mathbf{z}_{1:t-1}, \mathbf{Z}_t = k, \mathbf{x}_{1:t})$$

Gives what we want

$$\begin{aligned} & \max_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T}) \\ &= \max_{z_{1:T-1}, k} p(\mathbf{z}_{1:T-1}, \mathbf{Z}_T = k | \mathbf{x}_{1:T}) \\ &= \max_{z_{1:T-1}, k} \frac{p(\mathbf{z}_{1:T-1}, \mathbf{Z}_T = k, \mathbf{x}_{1:T})}{p(\mathbf{x}_{1:T})} \\ &= \max_{z_{1:T-1}, k} \frac{v_T(k)}{p(\mathbf{x}_{1:T})} \end{aligned}$$

VITERBI

- MAP path
- Viterbi learning: used, as approximation, to speed up parameter learning
- Again DP now with viterbi variable
- For the path, use back pointers

DP FOR VITERBI

$$\begin{aligned} v_t(k) &= \max_{z_{1:t-1}} p(\mathbf{z}_{1:t-1}, \mathbf{Z}_t = k, \mathbf{x}_{1:t}) \\ &= \max_l \max_{z_{1:t-2}} p(\mathbf{z}_{1:t-2}, \mathbf{Z}_{t-1} = l, \mathbf{Z}_t = k, \mathbf{x}_{1:t}) \\ &= \max_l \max_{z_{1:t-2}} p(\mathbf{z}_{1:t-2}, \mathbf{Z}_{t-1} = l, \mathbf{x}_{1:t-1}) p(\mathbf{Z}_t = k, \mathbf{x}_t | \mathbf{Z}_{t-1} = l) \\ &= \max_l \max_{z_{1:t-2}} \underbrace{v_{t-1}(l)}_{\text{smaller}} \underbrace{p(\mathbf{Z}_t = k | \mathbf{Z}_{t-1} = l)}_{\text{transition}} \underbrace{p(\mathbf{x}_t | \mathbf{Z}_t = k)}_{\text{emission}} \end{aligned}$$

- This gives the probability of the MAP path
- Using backpointers the path can be obtained

FORWARD FILTERING, BACKWARDS SAMPLING

$$\begin{aligned} & p(\mathbf{Z}_t = k | \mathbf{Z}_{t+1} = l, \mathbf{x}_{1:t}) \\ &= \frac{p(\mathbf{Z}_t = k, \mathbf{Z}_{t+1} = l, \mathbf{x}_{1:t})}{p(\mathbf{x}_{1:t}, \mathbf{Z}_{t+1} = l)} \\ \mathbf{z}_{1:T}^s \sim p(\mathbf{Z}_{1:T} | \mathbf{x}_{1:T}) &= \frac{p(\mathbf{x}_{1:t-1}, \mathbf{Z}_t = k) p(\mathbf{Z}_{t+1} = l, \mathbf{x}_t | \mathbf{Z}_t = k)}{f_{t+1}(l)} \\ &= \frac{f_t(k) p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) p(\mathbf{x}_t | \mathbf{Z}_t = k)}{f_{t+1}(l)} \end{aligned}$$

- Sample from posterior
- Sample in order z_T, \dots, z_1
- Start somewhat differently

LEARNING TRANSITION AND EMISSION PARAMETERS - FULLY OBSERVED DATA

- Parameters

- transition $A_{kl} = p(Z_t = l | Z_{t-1} = k)$
- z_0 always q^* and A_{q^*k} is start probability
- emission

$$B_{ks} = p(X_t = s | Z_t = k)$$

- Likelihood

$$L(\theta; \mathcal{D}) = \prod_{n=1}^N \prod_{t=1}^T \left[\prod_{k,s} B_{ks}^{I(x_{n,t}=s, z_{n,t}=k)} \prod_{k,l} A_{kl}^{I(z_{n,t-1}=k, z_{n,t}=l)} \right]$$

LEARNING TRANSITION AND EMISSION PARAMETERS - FULLY OBSERVED DATA

- Parameters

- transition $A_{kl} = p(Z_t = l | Z_{t-1} = k)$
- z_0 always q^* and A_{q^*k} is start probability
- emission

$$B_{ks} = p(X_t = s | Z_t = k)$$

- Loglikelihood

$$l(\theta; \mathcal{D}) = \sum_{n=1}^N \sum_{t=1}^T \left[\sum_{k,s} I(x_{n,t} = s, z_{n,t} = k) \log B_{ks} + \sum_{k,l} I(z_{n,t-1} = k, z_{n,t} = l) \log A_{kl} \right]$$

MAXIMIZING LOG-LIKELIHOOD - COMPLETE DATA

$$\begin{aligned} l(\theta; \mathcal{D}) &= \sum_{n=1}^N \sum_{t=1}^T \left[\sum_{k,s} I(x_{n,t} = s, z_{n,t} = k) \log B_{ks} + \sum_{k,l} I(z_{n,t-1} = k, z_{n,t} = l) \log A_{kl} \right] \\ &= \sum_{k,s} \left[\underbrace{\sum_{n=1}^N \sum_{t=1}^T I(x_{n,t} = s, z_{n,t} = k)}_{M_{k,s}} \right] \log B_{ks} + \sum_{k,l} \left[\underbrace{\sum_{n=1}^N \sum_{t=1}^T I(z_{n,t-1} = k, z_{n,t} = l)}_{N_{k,l}} \right] \log A_{kl} \end{aligned}$$

- Maximized by

$$B_{ks} = M_{k,s} / \sum_s M_{k,s} = M_{k,s} / N_k \quad \text{and} \quad A_{kl} = N_{k,l} / \sum_l N_{k,l} = N_{k,l} / N_k$$

HIDDEN VARIABLES - ONE EM-STEP MAXIMIZING Q

$$\begin{aligned} &\sum_{n=1}^N E_{p(\mathbf{Z} | \mathbf{x}_n, \theta)} [l(\theta'; \mathbf{Z}, \mathbf{x}_n)] \\ &= \sum_{n=1}^N E_{p(\mathbf{Z}, \mathbf{X} | \mathbf{x}_n, \theta)} [l(\theta'; \mathbf{Z}, \mathbf{X})] \\ &= \sum_{n=1}^N E_{p(\mathbf{Z}, \mathbf{X} | \mathbf{x}_n, \theta)} \left[\sum_{t=1}^T \left[\sum_{k,s} I(\mathbf{X}_t = s, \mathbf{Z}_t = k) \log B'_{ks} + \sum_{k,l} I(\mathbf{Z}_{t-1} = k, \mathbf{Z}_t = l) \log A'_{kl} \right] \right] \\ &= \sum_{k,s} \left[\underbrace{\sum_{n=1}^N \sum_{t=1}^T p(\mathbf{X}_t = s, \mathbf{Z}_t = k | \mathbf{X} = \mathbf{x}_n, \theta)}_{\bar{M}_{k,s}} \right] \log B'_{ks} \\ &\quad + \sum_{k,l} \left[\underbrace{\sum_{n=1}^N \sum_{t=1}^T p(\mathbf{Z}_{t-1} = k, \mathbf{Z}_t = l | \mathbf{X} = \mathbf{x}_n, \theta)}_{\bar{N}_{k,l}} \right] \log A'_{kl} \end{aligned}$$

LEARNING TRANSITION AND EMISSION PARAMETERS - HIDDEN VARIABLE

$$\begin{aligned} & \sum_{n=1}^N E_{p(\mathbf{Z}|\mathbf{x}_n, \theta)}[l(\theta'; \mathbf{Z}, \mathbf{x}_n)] \\ &= \sum_{k,s} \left[\underbrace{\sum_{n=1}^N \sum_{t=1}^T p(\mathbf{X}_t = s, \mathbf{Z}_t = k | \mathbf{X} = \mathbf{x}_n, \theta)}_{\bar{M}_{k,s}} \right] \log B'_{ks} \\ & \quad + \sum_{k,l} \left[\underbrace{\sum_{n=1}^N \sum_{t=1}^T p(\mathbf{Z}_{t-1} = k, \mathbf{Z}_t = l | \mathbf{X} = \mathbf{x}_n, \theta)}_{\bar{N}_{k,l}} \right] \log A'_{kl} \end{aligned} \quad \text{jl}$$

- Maximized by

$$B'_{ks} = M'_{k,s} / \sum_s M'_{k,s} = M'_{k,s} / N'_k \quad \text{and} \quad A'_{kl} = N'_{k,l} / \sum_l N'_{k,l} = N'_{k,l} / N'_k$$

THE END