

Bioinformatics
Introduction to Hidden Markov
Models

Hidden Markov Models and
Multiple Sequence Alignment

Slides borrowed from
Scott C. Schmidler
(MIS graduated student)

Outline

- Probability Review
- Markov Chains
- Hidden Markov Chains
- Examples in HMMs for Protein Sequence
- Algorithm Review for HMMs

Motivation: Composing a Drama by Mimicking Shakespeare

- Assume we want to write a drama of Shakespeare style
- We collect a large set of Shakespeare's works
- Define a vocabulary $V = \{X_1, X_2, \dots, X_N\}$
- Build a model $P(X_i|X_j)$ for $i, j = 1, \dots, N$
- To compose a drama, generate words from the model $P(X_i|X_j)$

- Though this is too simplistic to be useful, this naive model can be extended and refined to mimic the writing style of Shakespears'

Markov Approximations to English

- From Shannon's original paper:

1. *Zero-order approximation:*

XFOML RXKXRJFFUJ ZLPWCFWKCYJ FFJEYVKCQSGHYD
QPAAMKBZAACIBZLHJQD

2. *First-order approximation:*

OCRO HLI RGWR NWIELWIS EU LL NBNESEBYA TH EEI
ALHENHTTPA OOBTTVA NAH RBL

3. *Second-order approximation:*

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY
ACHIN D ILONASIVE TUCOOWE AT TEASONARE FUSO
TIZIN ANDY TOBE SEACE CITSBE

Markov Approximations (cont.)

From Shannon's paper

4. *Third-order approximation:*

IN NO IST LAT WHEY CRATICT FROURE BIRS GROCID
PONDENOME OF DEMONSTURES OF THE REPTABIN IS
REGOACTIONA OF CRE

Markov random field with 1000 "features," no underlying
"machine" (Della Pietra et. Al, 1997):

WAS REASER IN THERE TO WILL WAS BY HOMES THING
BE RELOVERATED THER WHICH CONISTS AT RORES
ANDITING WITH PROVERAL THE CHESTRAING FOR
HAVE TO INTRALLY OF QUT DIVERAL THIS OFFECT
INATEVER THIFER CONSTRADED STATER VILL
MENTTERING AND OF IN VERATE OF TO

Word-Based Approximations

1. *First-order approximation:*

REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME
CAN DIFFERENT NATURAL HERE HE THE A IN CAME THE TO
OF TO EXPERT GRAY COME TO FURNISHES THE LINE
MESSAGE HAD BE T

2. *Second-order approximation:*

THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH
WRITER THAT THE CHARACTER OF THIS POINT IS
THEREFORE ANOTHER METHOD FOR THE LETTERS THAT THE
TIME OF WHO EVER TOLD THE PROBLEM FOR AN
UNEXPETED

Shannon's comment:

“It would be interesting if further approximations could be constructed,
but the labor involved becomes enormous at the next stage.”

Motivation: Composing a Symphony of Beethoven Style

- We want to compose a symphony of Beethoven style
- We collect a large set of Beethoven's works
- Define a vocabulary $V = \{X_1, X_2, \dots, X_N\}$ of musical notes
- Build a model $P(X_i|X_j)$ for $i, j = 1, \dots, N$
- To compose a symphony, generate note symbols from the model $P(X_i|X_j)$

Modeling Biological Sequences

- Collect a set of sequences of interest
- Define a vocabulary $V = \{X_1, X_2, \dots, X_N\}$
 - ▶ For DNA sequences: $N = 4$ and $V = \{A, T, G, C\}$
 - ▶ For protein sequences: $N = 20$ and $V = \{\text{amino acids}\}$
- Build (learn) a model $P(X_i|X_j)$ for $i, j = 1, \dots, N$ or in more general $P(X|w)$ with $X = X_1, X_2, \dots, X_M$ and model parameter vector w
- The model can be used to
 - ▶ To generate typical sequences from the class of training sequences, e.g. protein family
 - ▶ To compute the probability of an observed sequence O being generated from the model class
 - ▶ and others
- Hidden Markov models (HMMs) are a class of stochastic generative models effective for building such probabilistic models.

Probability Review

- Probability notation:
 - ▶ Probability: $P(A) \quad P(A) \geq 0, \sum_A P(A) = 1$
 - ▶ Joint probability: $P(A, B)$
 - ▶ Conditional probability: $P(A | B) = P(A, B) / P(B)$
 - ▶ Marginal probability: $P(A) = \sum_B P(A, B)$
 - ▶ Independence: $P(A, B) = P(A)P(B)$
 - ▶ Bayes' rule: $P(B | A) = P(A | B)P(B) / P(A)$

Markov Chains

- Markov property:

$$P(X_0, X_1, \dots, X_t) = P(X_0)P(X_1 | X_0) \dots P(X_t | X_{t-1})$$

- Formally:

- ▶ State space

$$S = \{S_1, \dots, S_N\}$$

- ▶ Transition matrix

$$P = P(X_t = S_i | X_{t-1} = S_j) \quad 1 \leq i, j \leq N$$

- ▶ Initial distribution

$$\pi = P(X_0 = S_i) \quad 1 \leq i \leq N$$

- CS intuition

- ▶ Stochastic finite automaton



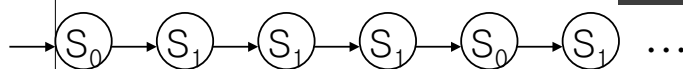
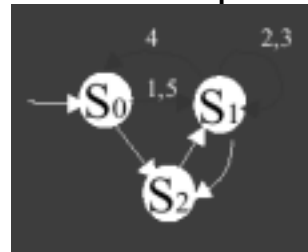
Markovian Sequence

- States through which the chain passes form a sequence:

Example: $S_0, S_1, S_1, S_1, S_0, S_1, \dots$

- Graphically:

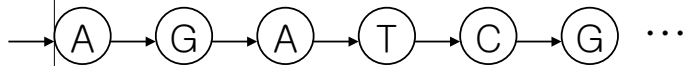
- By the Markov property:



$$\begin{aligned} P(\text{Sequence}) &= P(S_0, S_1, S_1, S_1, S_0, S_1, \dots) \\ &= \pi(S_0)P(S_1 | S_0)P(S_1 | S_1) \end{aligned}$$

Example

- Markov chain for generating a DNA sequence:
- Sequence probability:



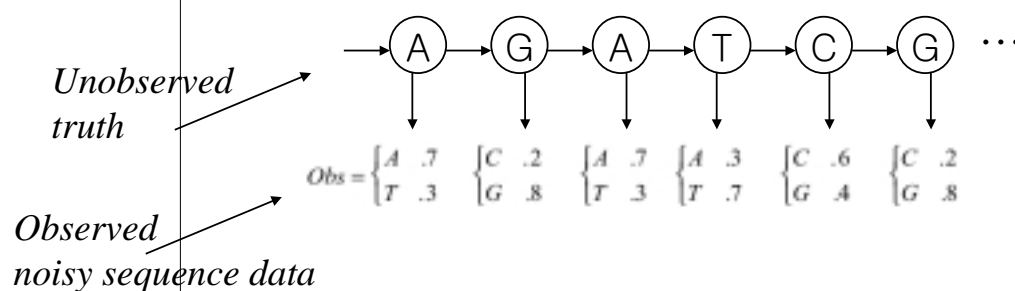
$$P(AGATCG) = \pi(A)P(G | A)P(A | G)P(T | A)P(C | T)P(G | C)$$

Dinucleotide frequency (e.g. base-stacking)

Hidden Markov Chains

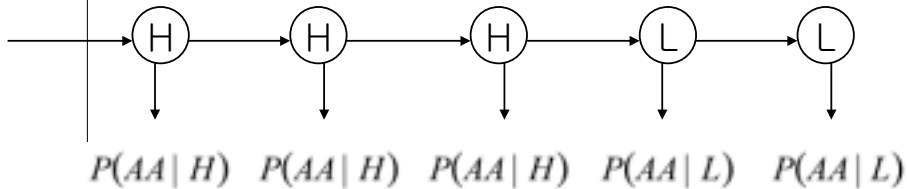
- Observed sequence is a *probabilistic function* of underlying Markov chain
 - ▶ Example: HMM for a (noisy) DNA sequence (see e.g. Churchill 1989)

True state sequence unknown, but observation sequence gives us a clue



Example: Hidden Markov Chain for Protein Sequence

- State space is backbone secondary structure
 - ▶ Used for prediction (Asai *et. al.*, Stultz *et. al.*)



- State space is side chain environment
 - ▶ Used for fold-recognition (Hubbard *et. al.*)

A HMM for Multiple Protein Sequences (Krogh *et. al.*)

- “Match” states are model (consensus) positions
- Position-specific deletion penalties
- Position-specific insertion frequencies
- Path through states aligns sequence to model

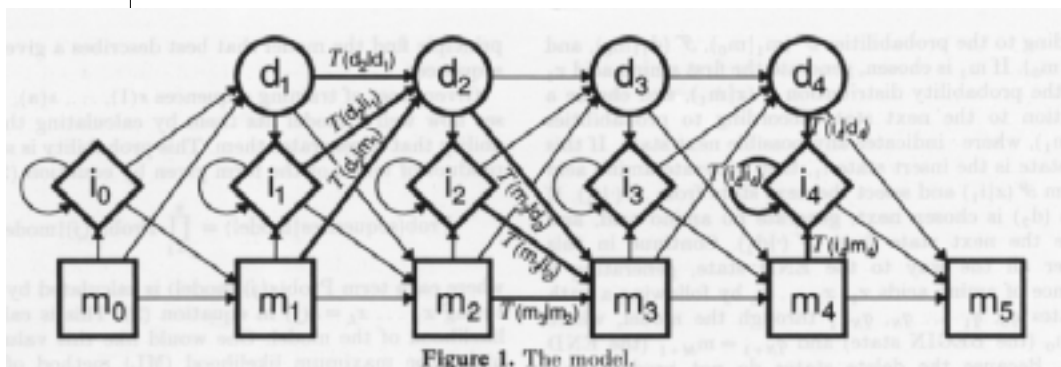


Figure from (Krogh *et. al.* 1994)

Example: Multiple Alignment of Globin Sequences

```

Helix      AAAAAAAAAAAAAAAAAA  BBBBBBBBBBBBBBBBCCCCCCCCCCCC  DDDDDDDDEE
*****
HBA_HUMAN  V.....LSPADKTNVKAANGKVG...NAGEYVGAEALGRMLSFPTTKTYFPHF-DLSHGSAQ----
HBB_HUMAN  Vh.....LTPEEKSAVTALWGKV--..NVDEVGGEALGRLLVVYPVTQRFFESFGDLSTPDAVHGSP
MYG_PHYCA  V.....LSEGEWQLVLRVWAKVEA...DVAGHGQDILIRLPKSHPETLEKFDNFKHLKTEAEMKASE
GLB3_CHITP  ....LSADQISTVQASFDK...KGDPSV--ILYAVPKADPSINAKFTQF-AGKDLESIKGTA
GLB5_PETMA  PivdtgsvapLSAAEKTAKIRSAWAPVYS..TYETSQVDILVKFPTSTPAAQEFPPKPKGLTTADQLKKA
LGB2_LUPLU  Ga.....LYESQAALVKSSWEFNA...NIPKNTHRFFILVLEIAPAADLF-SFLKGTSEVPQ--NNP
GLB1_GLYDI  G.....LSAAQNVIAATWRDIAGadngagVGKDCLIKFLSAEPQMAAVF-GF----SGASD---P

Helix      EEEEEEEEEEEEEEEEE  FFFFFFFF  FFFFFFFG  GGGGGGGGGGGGGGGG
*****
HBA_HUMAN  -VKGHGKKVADALTNAVAHVDD...MPNALSALSDLHA...HKLAVDPV..NFKLLSHCLLVTLAAHLP
HBB_HUMAN  KVKARHGKVLGAFSDGLAHLDN....LKGTFATLSELHC...DKLEVDPE..NPRLLGNVLVLCVLAHHPG
MYG_PHYCA  DLKKGVTVLTALGAILKKGH...HEAELKPLAQSHA...TK-EKIPiKYLEFISEAIINVLMSRNP
GLB3_CHITP  PFETHANRIVGFPSKIIIGELPN...IEADVNTFVASHK...PR-QVTHD..QLNNFRAGFVDSYMKHP
GLB5_PETMA  DVRWNAERIINAVNDAVASMDDeek..NSMKLSDLSGKHA...KSPQVDPQ..YFKVLAAVIADTVAA---
LGB2_LUPLU  ELQAHAGKVFKLVEAAIQLOVtgvvvTDA TLKNLGSHV...SK-GVADA..HFPVVEAAILKTIKEVVG
GLB1_GLYDI  GVAALGAKVLAQIGVAVSHLGDegk..MVAQNKAVGVRHNggNK-EIKAQ..YFEPLGASLLSMEHRRIG

Helix      HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH
*****
HBA_HUMAN  AEFTPAVHASLDRFLASVSTVLTISKY.....R
HBB_HUMAN  KEFTPPVQAAYQKVVAGVANALANKY.....H
MYG_PHYCA  GDFGADAQGAMNKALELFRKDIAAKYkelgyqG
GLB3_CHITP  TDF-AGAEAANGATLDTFFGMIFSKM.....-
GLB5_PETMA  GD-----AGFERLNSMICILLRSAY.....-
LGB2_LUPLU  AKWSEKLSAWTIAYDELAIVIKKEMnda....A
GLB1_GLYDI  GKMNAAAKDAVAAAYADISGALISGLq....S
    
```

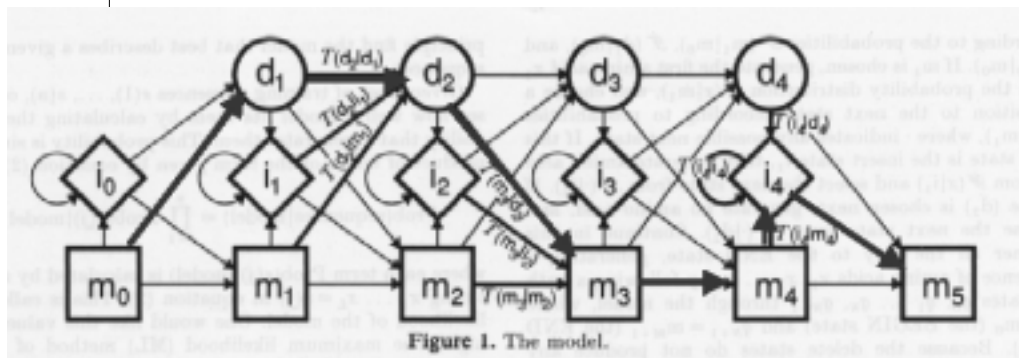
Figure from (Krogh *et al.* 1994)

HMM-based Multiple Sequence Alignment

- Multiple alignment of k sequences is $O(n^k)$, so instead:
 1. Estimate a statistical model for the sequences
 - Use head start PROFILE alignment
 - Start from scratch with unaligned sequences (harder)
 2. Align each remaining sequence to the model
 3. Alignment yields assignments of equivalent sequence elements within the multiple alignment

Example: Aligning Sequence to Model

- Given an HMM model for a protein family:
Align a new sequence to the model
(d states are gaps, i states are insertions)



(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

19

Computing with HMMs

- Three tasks:

1. Probability of an observed sequence

Given O_1, O_2, \dots, O_r find $P(O_1, O_2, \dots, O_r)$

(nontrivial since state sequence unobserved)

2. Most likely hidden state sequence

Given O_1, O_2, \dots, O_r compute $\arg \max_{S_1, K, S_r} P(S_1, K, S_r | O_1, O_2, \dots, O_r)$

2. Most likely hidden state sequence

Given observed sequence $\{O^1, \dots, O^n\}$ find

$$\arg \max_{\theta} P(O^1, K, O^n | \theta)$$

(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

20

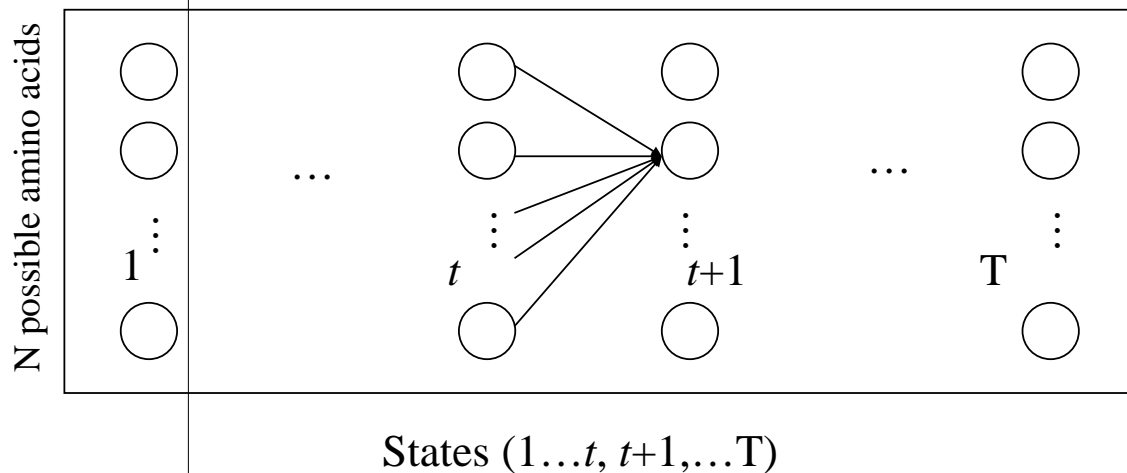
Computing Likelihood of Observed Sequence

- Compute $P(O_1, O_2, \dots, O_r)$
 - ▶ True state sequence unknown
 - ▶ Must sum over all possible paths
 - ▶ Number paths $O(T^N)$
 - ▶ Markovian structure permits:
 - Recursive definition and hence
 - Efficient calculation by dynamic programming

$$P(O_1, \dots, O_r) = \sum_{S_0, S_1, \dots, S_r} P(O_1, O_2, \dots, O_r | S_0, S_1, \dots, S_r) P(S_0, S_1, \dots, S_r)$$

- Key observation: Any path must be in exactly one state at time t

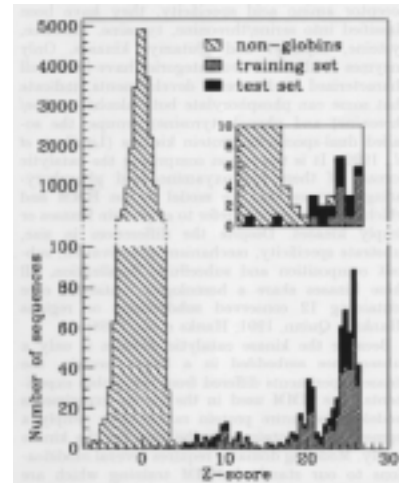
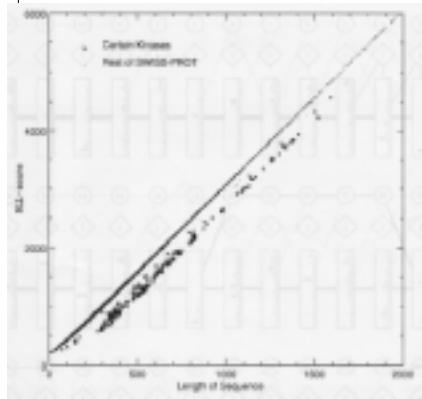
Key Idea for HMM Computations



Example: Searching Protein Database with HMM Profile

- For each sequence in database:
- Does sequence “fit” model?
- Score by $P(O_1, O_2, \dots, O_r)$, compute Z-score adjusted for length

Protein Kinases:



Globins:

Figure from (Krogh *et al.* 1994)

(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

23

Estimate Alignment and Model Parameters - *Simultaneously*

- Key idea – missing data:
 - ▶ *What if we know the alignment?*
 - Parameters easy to estimate:
 - Calculate (expected) number of transitions
 - Calculate (expected) frequency of amino acids
 - ▶ *What if we knew the parameters?*
 - Alignment easy to find
 - Align each sequence to model using Viterbi algorithm
 - Align residues in match states

(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

24

Other “details”

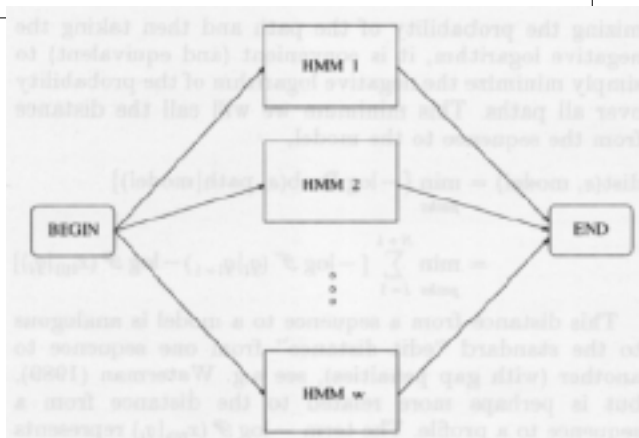
- How many states in model?
- How to initialize parameters?
- How to avoid local models?

See (Krong *et. al.*, 1994) for some suggestion

Multiple Protein Sequence Alignment

- Give a set of sequences:
 - ▶ Estimate HMM model using optimization for parameter search (Baum-Welch, EM)
 - ▶ Align each sequence to model (Viterbi)
 - ▶ Match states of model provide columns of resulting multiple alignment

Extensions



Clustering subfamilies

Modeling domains

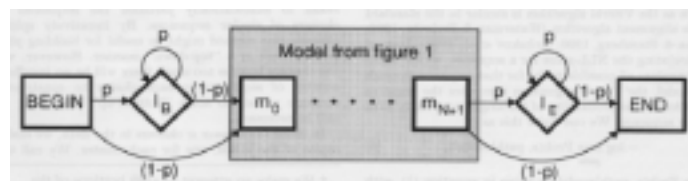


Figure from (Krogh *et. al.* 1994)

(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

27

Tradeoffs

- Advantages:
 - ▶ Explicit probabilistic model for family
 - ▶ Position specific residue distributions, gap penalties, insertions frequencies
- Disadvantages:
 - ▶ Many parameters, requires more data of care
 - ▶ Traded one hard optimization problem for another

(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

28

HMM Summary

- Powerful tool for modeling protein families
- Generalization of existing profile methods
- Data-intensive
- Widely applicable to problems in bioinformatics

References

- **Bioinformatics Classic: Krogh *et. al.* (1994)**
Hidden Markov models in computational biology: applications to protein modeling, *J. Mol. Biol.* 235: 1501-1531
- **Book: Eddy & Durbin, 1999. See web site.**
- **Tutorial: Rabiner, L. (1989) A tutorial on hidden Markov models and selected applications in speech recognition,**
Proc IEEE, 77(2), 257-286

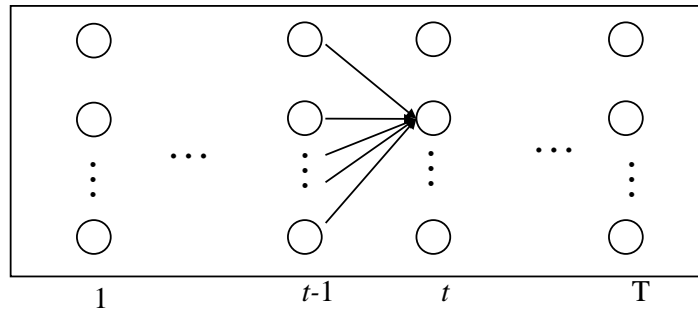
Forward-backward Algorithm

- Forward pass:

- ▶ Define $\alpha_t(j) = \left[\sum_{i=1}^N \alpha_{t-1}(i) P(S_j | S_i) \right] P(O_t | S_j)$

- ▶ Prob. Of subsequence O_1, O_2, \dots, O_t when in S_j at t

Key obs: any path must be in 1 of N states at t



(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

31

Forward-backward Algorithm

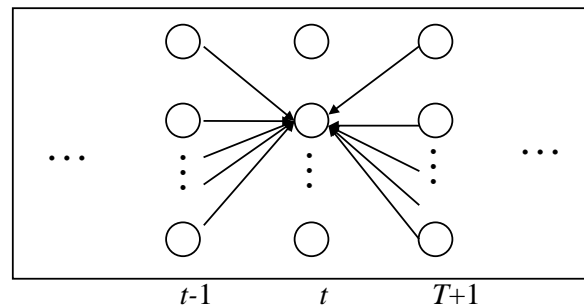
- Notice $P(O_1, O_2, \dots, O_r) = \sum_{j=1}^N \alpha_r(j)$

- Define an analogous backward pass so that:

$$\beta_t(j) = \sum_{i=1}^N \beta_{t+1}(i) P(S_i | S_j) P(O_{t+1} | S_i)$$

and

$$P(O_t \text{ came from } S_j) = \frac{\alpha_t(j) \beta_t(j)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$



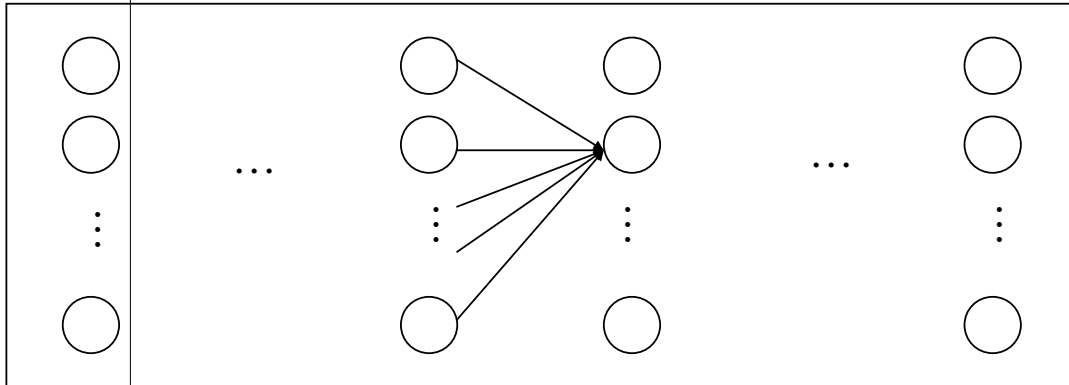
(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

32

Finding Most Likely Path

- Forward pass:

- ▶ Replace summation with maximization
- ▶ Max prob. of subseq. O_1, O_2, \dots, O_r When in S_j at t
- ▶ Again: $\max P(O_1, O_2, \dots, O_r) = \max_{1 \leq j \leq N} \alpha_T(j)$, then trace back



(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

33

Baum-Welch Algorithm (Expectation-Maximization)

- Set parameters to expected values given observed sequences:

- ▶ State transition probs:

- ▶ Observation probs:

$$P(S_j | S_i) = \frac{\sum_{t=1}^{T-1} P(\text{in } S_i \text{ at } t, \text{in } S_j \text{ at } t+1 | O)}{\sum_{t=1}^{T-1} P(\text{in } S_i \text{ at } t | O)}$$

$$P(\text{obs} | S_j) = \frac{\sum_{t=1}^{T-1} P(\text{in } S_j \text{ at } t | O) * 1(O_t = \text{obs})}{\sum_{t=1}^{T-1} P(\text{in } S_j \text{ at } t | O)}$$

- ▶ Recalculate expectations with new probabilities

- ▶ Iterate to convergence

- Guaranteed $P(O^1, K, O^n | \theta)$ strictly increasing, converge to local mode (See Rabiner, 1989 for details)

(C) 2001 SNU CSE Artificial Intelligence Lab (SCAI)

34