

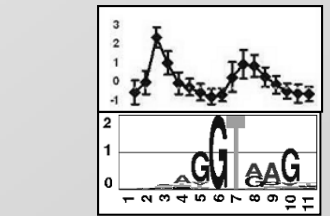
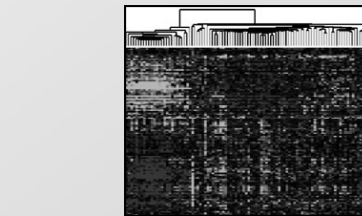
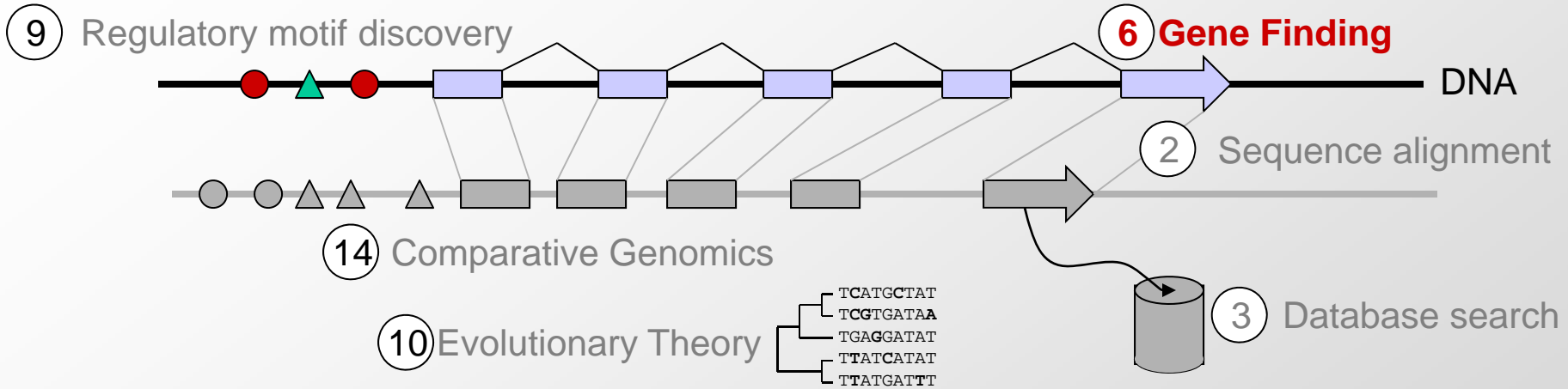
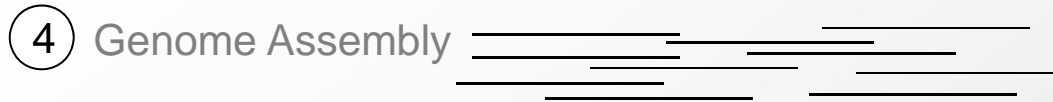
MIT OpenCourseWare
<http://ocw.mit.edu>

6.047 / 6.878 Computational Biology: Genomes, Networks, Evolution
Fall 2008

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.

Modeling Biological Sequence using Hidden Markov Models

Challenges in Computational Biology



⑫ Regulatory network inference

⑬ Emerging network properties

What have we learned so far?

- String searching and counting
 - Brute-force algorithm
 - W-mer indexing
- Sequence alignment
 - Dynamic programming, duality path \Leftrightarrow alignment
 - Global / local alignment, general gap penalties
- String comparison
 - Exact string match, semi-numerical matching
- Rapid database search
 - Exact matching: Hashing, BLAST
 - Inexact matching: neighborhood search, projections
- Problem set 1

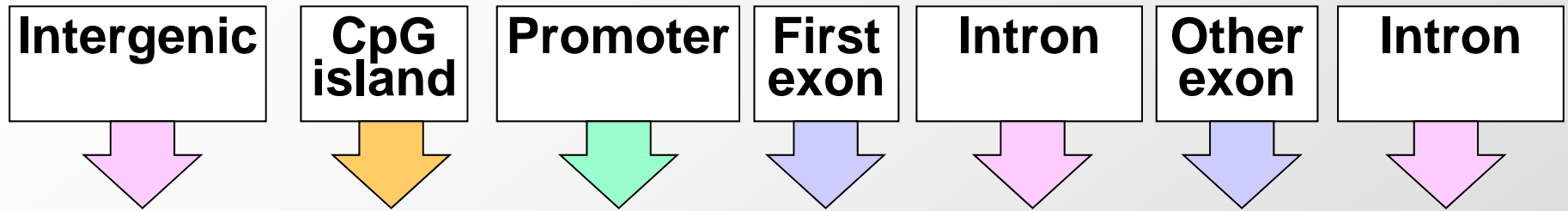
So, you find a new piece of DNA...

What do you do?

...GTACTCACCGGGTTACAGGATTATGGGTTACAGGTAACCGTT...

- Align it to things we know about
- Align it to things we don't know about
- Stare at it
 - Non-standard nucleotide composition?
 - Interesting k-mer frequencies?
 - Recurring patterns?
- Model it
 - Make some hypotheses about it
 - Build a 'generative model' to describe it
 - Find sequences of similar *type*

This week: Modeling biological sequences (a.k.a. What to do with a huge chunk of DNA)



- Ability to **emit** DNA sequences of a certain *type*
 - Not exact alignment to previously known gene
 - Preserving ‘properties’ of **type**, not identical sequence
- Ability to **recognize** DNA sequences of a certain type (state)
 - What (hidden) state is most likely to have generated observations
 - Find set of states and transitions that generated a long sequence
- Ability to **learn** distinguishing characteristics of each state
 - Training our generative models on large datasets
 - Learn to classify unlabelled data

Why Probabilistic Sequence Modeling?

- Biological data is noisy
- Probability provides a calculus for manipulating models
- Not limited to yes/no answers – can provide “degrees of belief”
- Many common computational tools based on probabilistic models
- Our tools:
 - Markov Chains and Hidden Markov Models (HMMs)

Definition: Markov Chain

Definition: A *Markov chain* is a triplet $(\mathbf{Q}, \mathbf{p}, \mathbf{A})$, where:

- \mathbf{Q} is a finite set of states. Each state corresponds to a symbol in the alphabet Σ
- \mathbf{p} is the initial state probabilities.
- \mathbf{A} is the state transition probabilities, denoted by \mathbf{a}_{st} for each \mathbf{s}, \mathbf{t} in \mathbf{Q} .
- For each \mathbf{s}, \mathbf{t} in \mathbf{Q} the transition probability is: $\mathbf{a}_{st} \equiv P(\mathbf{x}_i = \mathbf{t} | \mathbf{x}_{i-1} = \mathbf{s})$

Output: The output of the model is the set of states at each instant time => the set of states are observable

Property: The probability of each symbol \mathbf{x}_i depends only on the value of the preceding symbol \mathbf{x}_{i-1} : $P(\mathbf{x}_i | \mathbf{x}_{i-1}, \dots, \mathbf{x}_1) = P(\mathbf{x}_i | \mathbf{x}_{i-1})$

Formula: The probability of the sequence:

$$P(\mathbf{x}) = P(\mathbf{x}_L, \mathbf{x}_{L-1}, \dots, \mathbf{x}_1) = P(\mathbf{x}_L | \mathbf{x}_{L-1}) P(\mathbf{x}_{L-1} | \mathbf{x}_{L-2}) \dots P(\mathbf{x}_2 | \mathbf{x}_1) P(\mathbf{x}_1)$$

Definitions: HMM (Hidden Markov Model)

Definition: An *HMM* is a 5-tuple (Q, V, p, A, E) , where:

- Q is a finite set of states, $|Q|=N$
- V is a finite set of observation symbols per state, $|V|=M$
- p is the initial state probabilities.
- A is the state transition probabilities, denoted by a_{st} for each s, t in Q .
 - For each s, t in Q the transition probability is: $a_{st} \equiv P(x_i = t | x_{i-1} = s)$
- E is a probability emission matrix, $e_{sk} \equiv P(v_k \text{ at time } t | q_t = s)$

Output: Only **emitted symbols** are observable by the system but not the underlying random walk between states -> "hidden"

Property: **Emissions** and **transitions** are dependent on the current state only and not on the past.

The six algorithmic settings for HMMs

One path

All paths

Scoring

1. Scoring x , one path

$$P(x, \pi)$$

Prob of a path, emissions

2. Scoring x , all paths

$$P(x) = \sum_{\pi} P(x, \pi)$$

Prob of emissions, over all paths

Decoding

3. Viterbi decoding

$$\pi^* = \operatorname{argmax}_{\pi} P(x, \pi)$$

Most likely path

4. Posterior decoding

$$\pi^{\wedge} = \{ \pi_i \mid \pi_i = \operatorname{argmax}_k \sum_{\pi} P(\pi_i = k | x) \}$$

Path containing the most likely state at any time point.

Learning

5. Supervised learning, given π

$$\Lambda^* = \operatorname{argmax}_{\Lambda} P(x, \pi | \Lambda)$$

6. Unsupervised learning.

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \max_{\pi} P(x, \pi | \Lambda)$$

Viterbi training, best path

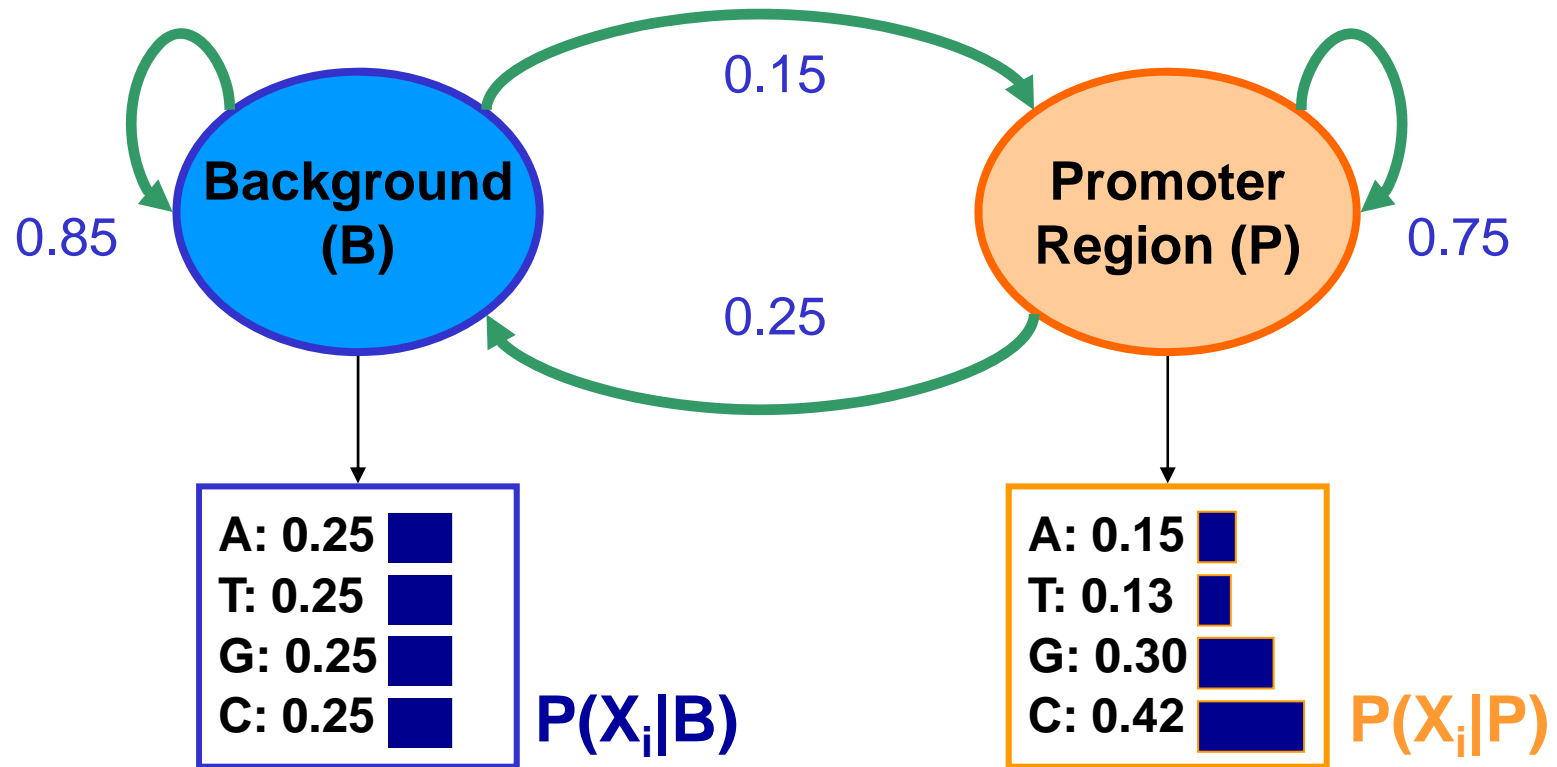
6. Unsupervised learning

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \sum_{\pi} P(x, \pi | \Lambda)$$

Baum-Welch training, over all paths

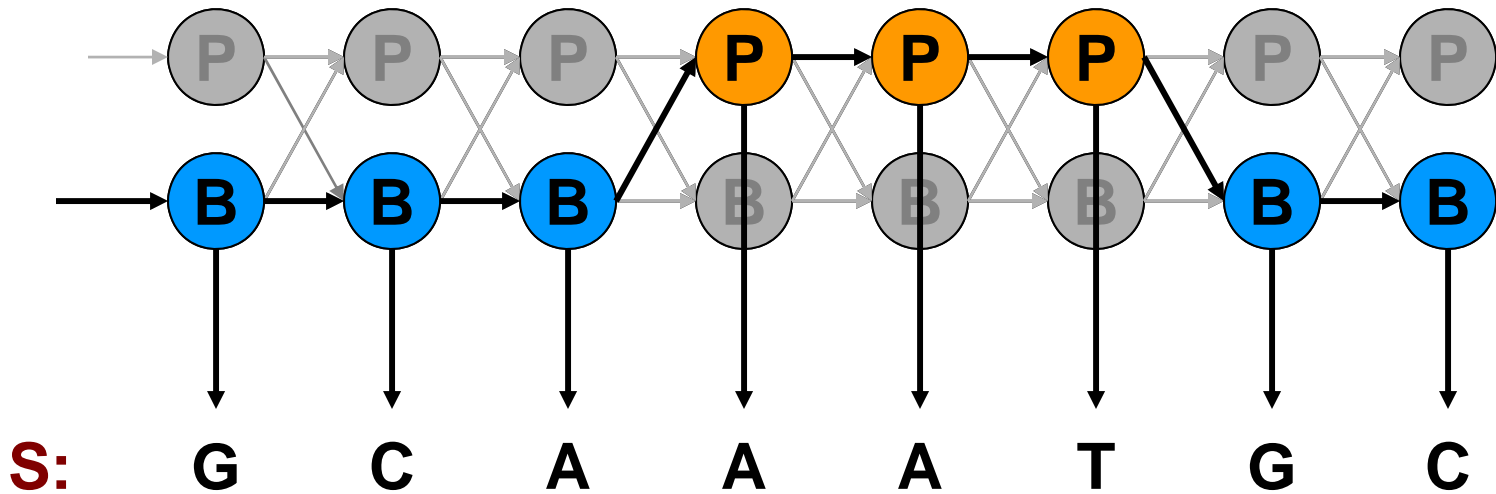
Example 1: Finding GC-rich regions

- Promoter regions frequently have higher counts of Gs and Cs
- Model genome as nucleotides drawn independently from two distributions: Background (B) and Promoters (P).
- Emission probabilities based on nucleotide composition in each.
- Transition probabilities based on relative abundance & avg. length



TAAGAATTGTGTCACACACATA**AAAACCCT**AAGTTAGAGGATTGAGATTGGCA
GACGATTGTT**CGTGATAATAAAC**AAGGGGGGCATAGATCAGGCTCATATTGGC

HMM as a *Generative* Model



$P(L_{i+1}|L_i)$

	B_{i+1}	P_{i+1}
B_i	0.85	0.15
P_i	0.25	0.75

$P(S|B)$

A:	0.25	■
T:	0.25	■
G:	0.25	■
C:	0.25	■

$P(S|P)$

A:	0.42	■
T:	0.30	■
G:	0.13	■
C:	0.15	■

Sequence Classification

PROBLEM: Given a sequence, is it a promoter region?

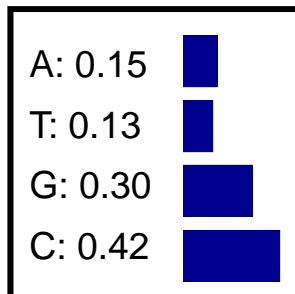
- We can calculate $P(S|MP)$, but what is a *sufficient P value*?

SOLUTION: compare to a **null model** and calculate **log-likelihood ratio**

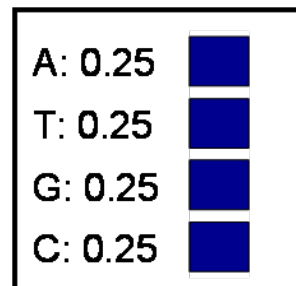
- e.g. background DNA distribution model, B

$$Score = \log \frac{P(S | MP)}{P(S | B)}$$

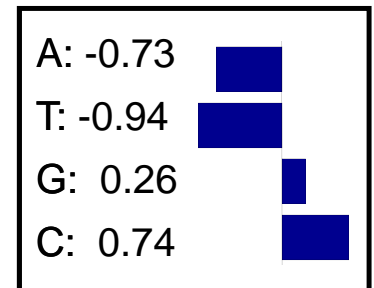
**Pathogenicity
Islands**



**Background
DNA**



**Score
Matrix**

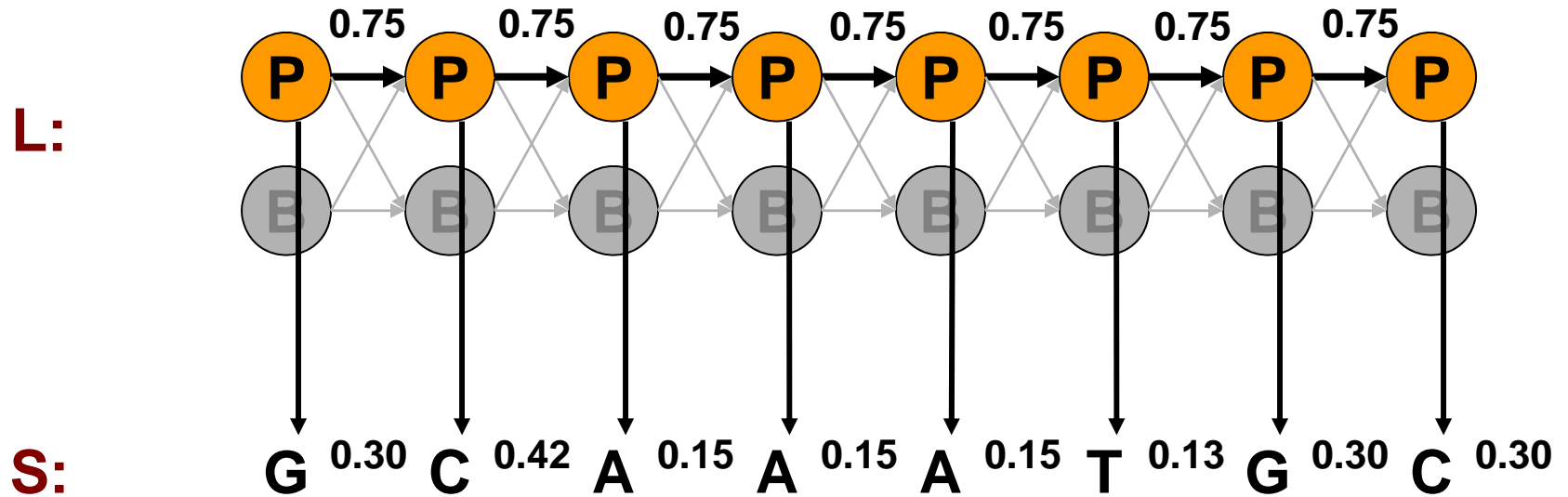


Finding GC-rich regions

- Could use the log-likelihood ratio on windows of fixed size
- Downside: have to evaluate all islands of all lengths repeatedly
- Need: a way to easily find transitions

TAAGAATTGTGTCACACACATAAAAACCTAAGTTAGAGGATTGAGATTGGCA
GACGATTGTCGTGATAATAACAAGGGGGGCATAGATCAGGCTCATATTGGC

Probability of a sequence if all promoter

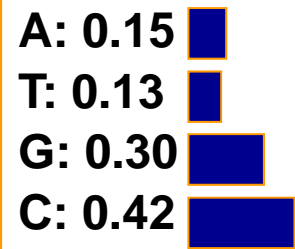


$$P(x, \pi) = a_p * e_p(G) * a_{pp} * e_p(G) * a_{pp} * e_p(C) * a_{pp} * e_p(A) * a_{pp} * \dots$$

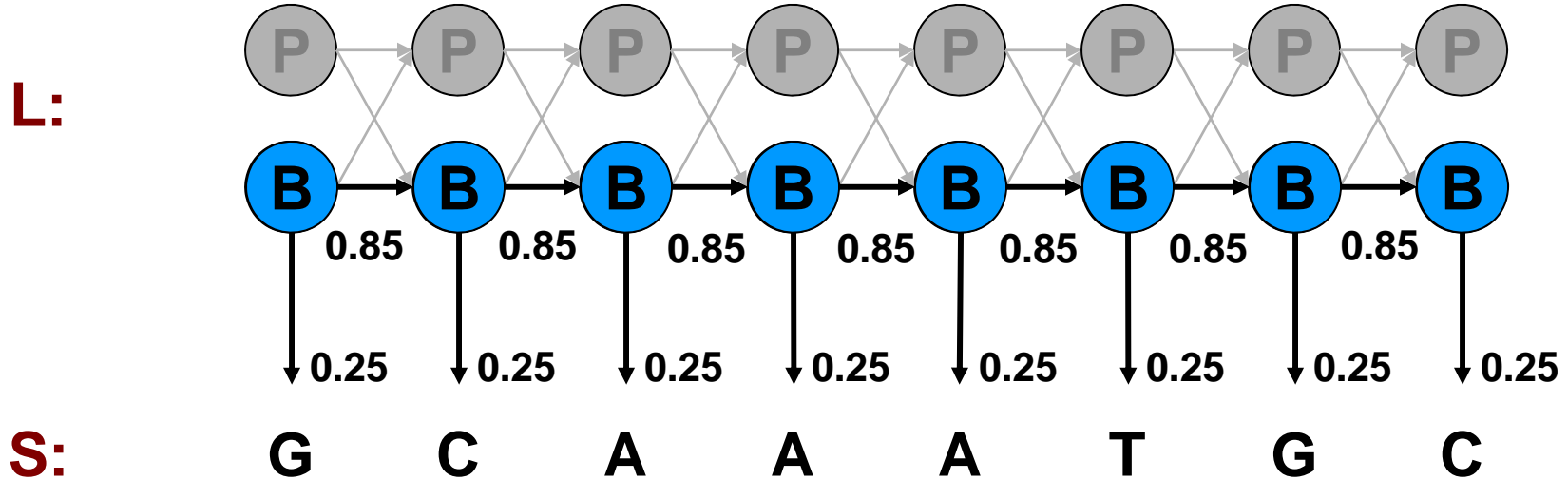
$$= a_p * (0.75)^7 * (0.15)^3 * (0.13)^1 * (0.30)^2 * (0.42)^2$$

$$= 9.3 * 10^{-7}$$

Why is this so small?



Probability of the same sequence if all background

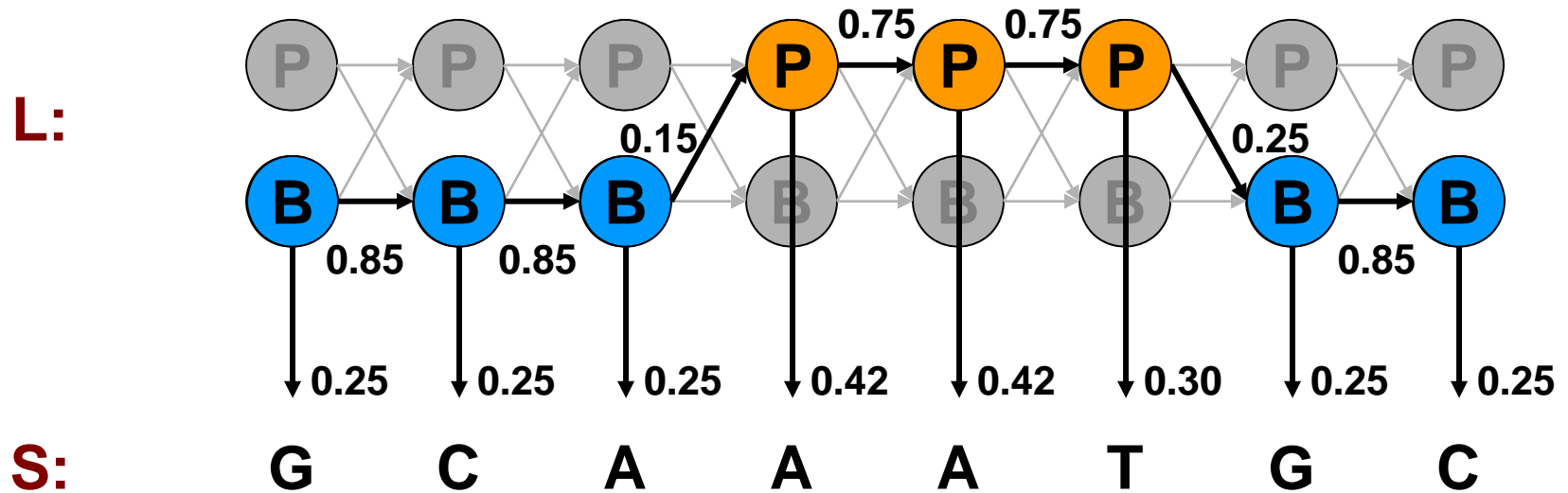


$$\begin{aligned} P &= P(G | B)P(B_1 | B_0)P(C | B)P(B_2 | B_1)P(A | B)P(B_3 | B_2)\dots P(C | B_7) \\ &= (0.85)^7 \times (0.25)^8 \\ &= 4.9 \times 10^{-6} \end{aligned}$$

Compare relative probabilities: 5-fold more likely!

A: 0.25	■
T: 0.25	■
G: 0.25	■
C: 0.25	■

Probability of the same sequence if mixed



$$\begin{aligned}
 P &= P(G | B)P(B_1 | B_0)P(C | B)P(B_2 | B_1)P(A | B)P(P_3 | B_2)...P(C | B_7) \\
 &= (0.85)^3 \times (0.25)^6 \times (0.75)^2 \times (0.42)^2 \times 0.30 \times 0.15 \\
 &= 6.7 \times 10^{-7}
 \end{aligned}$$

Should we try all possibilities? What is the most likely path?

The six algorithmic settings for HMMs

One path

All paths

Scoring

1. Scoring x , one path

$$P(x, \pi)$$

Prob of a path, emissions

2. Scoring x , all paths

$$P(x) = \sum_{\pi} P(x, \pi)$$

Prob of emissions, over all paths

Decoding

3. Viterbi decoding

$$\pi^* = \operatorname{argmax}_{\pi} P(x, \pi)$$

Most likely path

4. Posterior decoding

$$\pi^{\wedge} = \{ \pi_i \mid \pi_i = \operatorname{argmax}_k \sum_{\pi} P(\pi_i = k | x) \}$$

Path containing the most likely state at any time point.

Learning

5. Supervised learning, given π

$$\Lambda^* = \operatorname{argmax}_{\Lambda} P(x, \pi | \Lambda)$$

6. Unsupervised learning.

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \max_{\pi} P(x, \pi | \Lambda)$$

Viterbi training, best path

6. Unsupervised learning

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \sum_{\pi} P(x, \pi | \Lambda)$$

Baum-Welch training, over all paths

3. DECODING:

What was the sequence of hidden states?

Given: Model parameters $e_i(\cdot)$, a_{ij}

Given: Sequence of emissions x

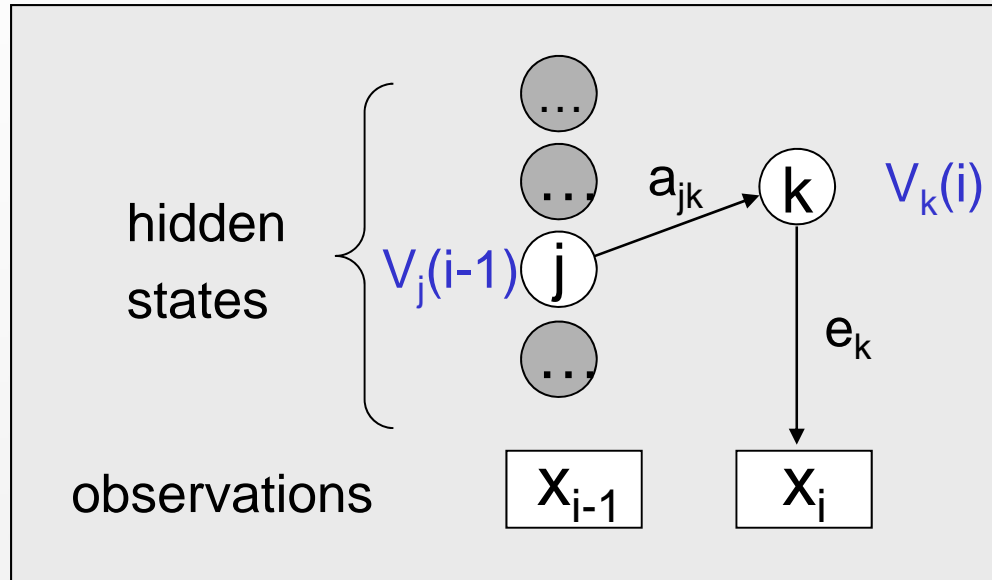
Find: Sequence of hidden states π

Finding the optimal path

- We can now evaluate any path through hidden states, given the emitted sequences
- How do we find the best path?
- Optimal substructure! Best path through a given state is:
 - Best path to previous state
 - Best transition from previous state to this state
 - Best path to the end state
- Viterbi algorithm
 - Define $V_k(i)$ = Probability of the most likely path through state $\pi_i=k$
 - Compute $V_k(i+1)$ as a function of $\max_{k'} \{ V_{k'}(i) \}$
 - $V_k(i+1) = e_k(x_{i+1}) * \max_j a_{jk} V_j(i)$

→ Dynamic Programming

Calculate maximum $P(x, \pi)$ recursively



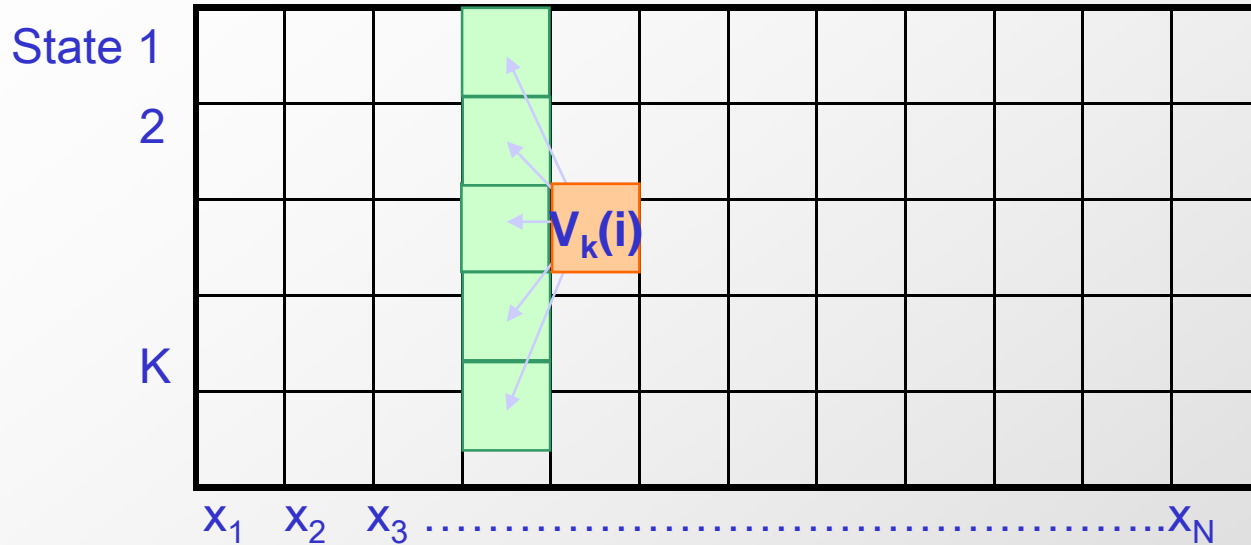
- Assume we know V_j for the previous time step $(i-1)$

- Calculate $V_k(i) = e_k(x_i) * \max_j (V_j(i-1) \times a_{jk})$

current max this emission max ending in state j at step i Transition from state j

all possible previous states j

The Viterbi Algorithm



Input: $x = x_1 \dots x_N$

Initialization:

$$V_0(0)=1, V_k(0) = 0, \text{ for all } k > 0$$

Iteration:

$$V_k(i) = e_K(x_i) \times \max_j a_{jk} V_j(i-1)$$

Termination:

$$P(x, \pi^*) = \max_k V_k(N)$$

Traceback:

Follow max pointers back

Similar to aligning states to seq

In practice:

Use log scores for computation

Running time and space:

Time: $O(K^2N)$

Space: $O(KN)$

The six algorithmic settings for HMMs

One path

All paths

Scoring

1. Scoring x , one path

$$P(x, \pi)$$

Prob of a path, emissions

2. Scoring x , all paths

$$P(x) = \sum_{\pi} P(x, \pi)$$

Prob of emissions, over all paths

Decoding

3. Viterbi decoding

$$\pi^* = \operatorname{argmax}_{\pi} P(x, \pi)$$

Most likely path

4. Posterior decoding

$$\pi^{\wedge} = \{ \pi_i \mid \pi_i = \operatorname{argmax}_k \sum_{\pi} P(\pi_i = k | x) \}$$

Path containing the most likely state at any time point.

Learning

5. Supervised learning, given π

$$\Lambda^* = \operatorname{argmax}_{\Lambda} P(x, \pi | \Lambda)$$

6. Unsupervised learning.

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \max_{\pi} P(x, \pi | \Lambda)$$

Viterbi training, best path

6. Unsupervised learning

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \sum_{\pi} P(x, \pi | \Lambda)$$

Baum-Welch training, over all paths

2. EVALUATION

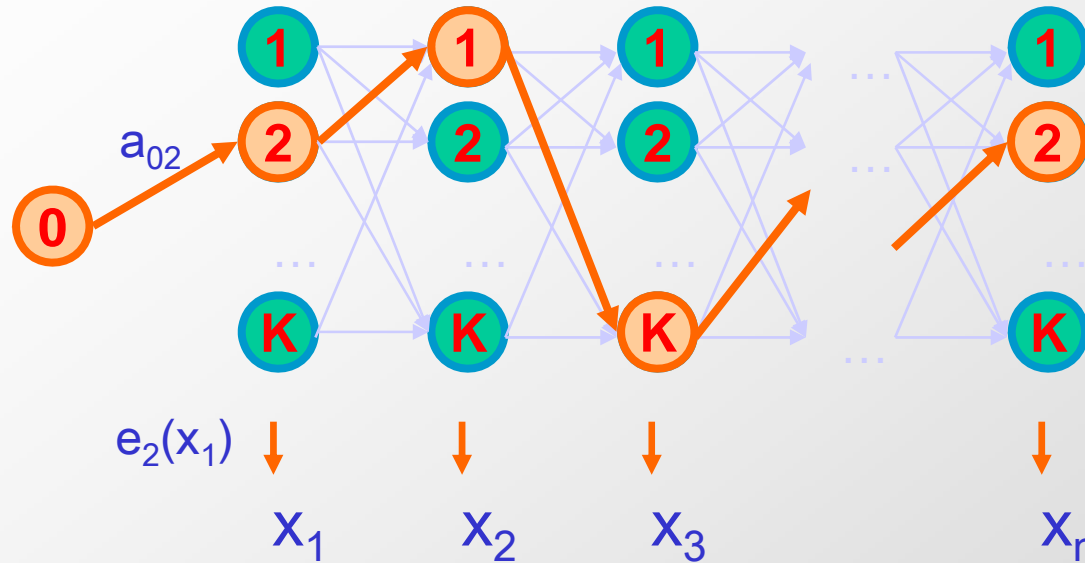
(how well does our model capture the world)

Given: Model parameters $e_i(\cdot)$, a_{ij}

Given: Sequence of emissions x

Find: $P(x|M)$, summed over all possible paths π

Simple: Given the model, generate some sequence x

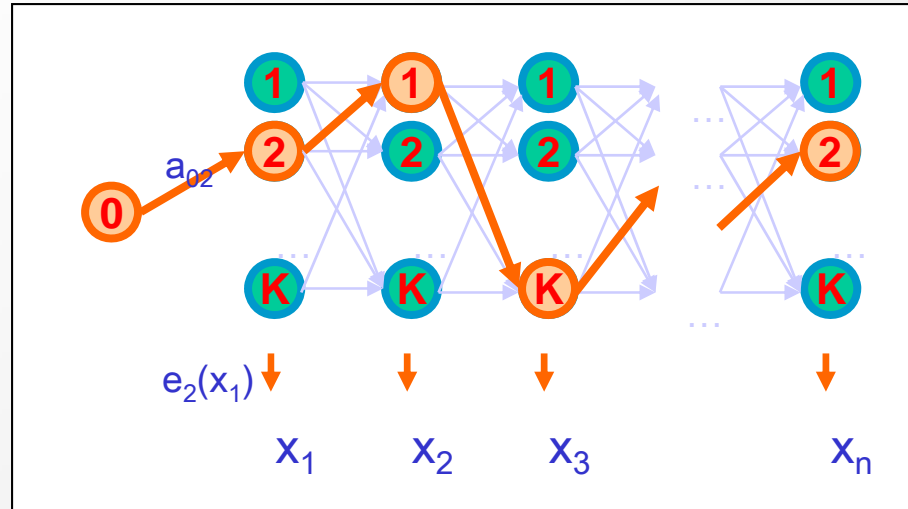


Given a HMM, we can generate a sequence of length n as follows:

1. Start at state π_1 according to prob $a_{0\pi_1}$
2. Emit letter x_1 according to prob $e_{\pi_1}(x_1)$
3. Go to state π_2 according to prob $a_{\pi_1\pi_2}$
4. ... until emitting x_n

We have some sequence x that can be emitted by p . Can calculate its likelihood. However, in general, many different paths may emit this same sequence x . How do we find the total probability of generating a given x , over any path?

Complex: Given x , was it generated by the model?



Given a sequence x ,

What is the probability that x was generated by the model (using any path)?

$$- P(x) = \sum_{\pi} P(x, \pi)$$

- Challenge: exponential number of paths

Calculate probability of emission over all paths

- Each path has associated probability
 - Some paths are likely, others unlikely: sum them all up
 - Return total probability that emissions are observed, summed over all paths
 - Viterbi path is the most likely one
 - How much 'probability mass' does it contain?
- (cheap) alternative:
 - Calculate probability over maximum (Viterbi) path π^*
 - Good approximation if Viterbi has highest density
 - BUT: incorrect
- (real) solution
 - Calculate the exact sum iteratively
 - $P(x) = \sum_{\pi} P(x, \pi)$
 - Can use dynamic programming

The Forward Algorithm – derivation

Define the forward probability:

$$f_l(i) = P(x_1 \dots x_i, \pi_i = l)$$

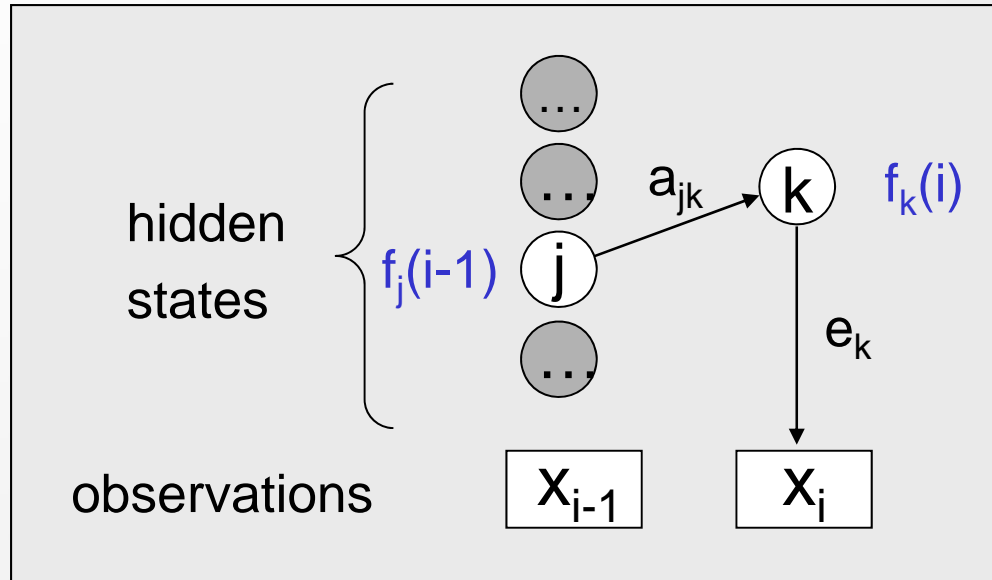
$$= \sum_{\pi_1 \dots \pi_{i-1}} P(x_1 \dots x_{i-1}, \pi_1, \dots, \pi_{i-2}, \pi_{i-1}, \pi_i = l) e_l(x_i)$$

$$= \sum_k \boxed{\sum_{\pi_1 \dots \pi_{i-2}} P(x_1 \dots x_{i-1}, \pi_1, \dots, \pi_{i-2}, \pi_{i-1} = k)} a_{kl} e_l(x_i)$$

$$= \sum_k \boxed{f_k(i-1)} a_{kl} e_l(x_i)$$

$$= e_l(x_i) \sum_k \boxed{f_k(i-1)} a_{kl}$$

Calculate total probability $\sum_{\pi} P(x, \pi)$ recursively



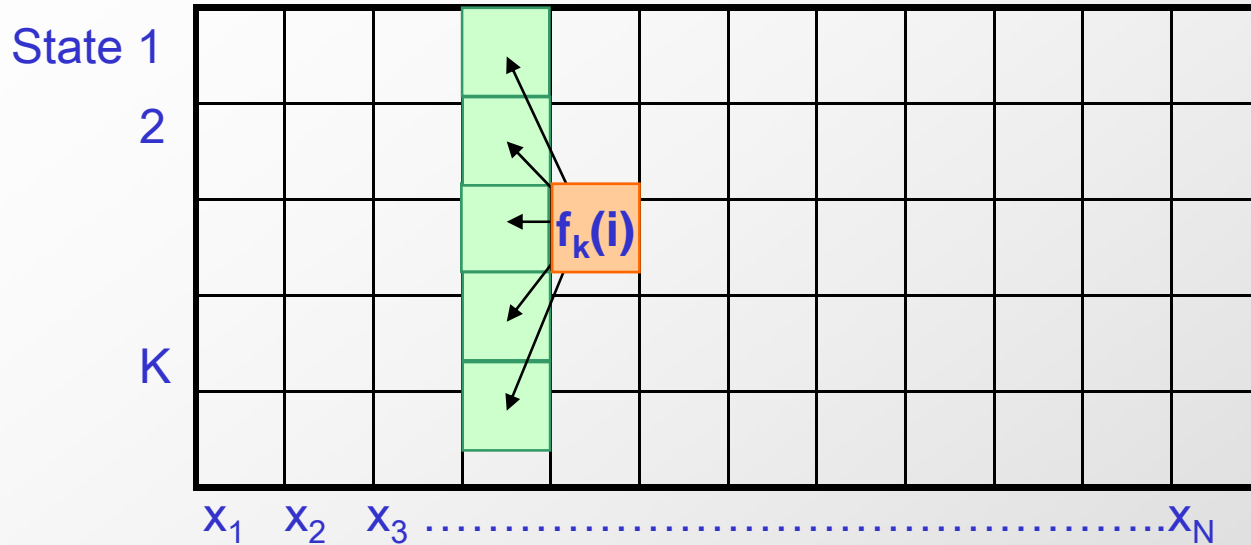
- Assume we know f_j for the previous time step $(i-1)$

- Calculate $f_k(i) = e_k(x_i) * \sum_j (f_j(i-1) \times a_{jk})$

updated sum
this emission
sum ending in state j at step i
transition from state j

every possible previous state j

The Forward Algorithm



Input: $x = x_1 \dots x_N$

Initialization:

$$f_0(0) = 1, f_k(0) = 0, \text{ for all } k > 0$$

Iteration:

$$f_k(i) = e_K(x_i) \times \sum_j a_{jk} f_j(i-1)$$

Termination:

$$P(x, \pi^*) = \sum_k f_k(N)$$

In practice:

- Sum of log scores is difficult
- approximate $\exp(1+p+q)$
- scaling of probabilities

Running time and space:

Time: $O(K^2N)$

Space: $O(KN)$

The six algorithmic settings for HMMs

One path

All paths

Scoring

1. Scoring x , one path

$$P(x, \pi)$$

Prob of a path, emissions

2. Scoring x , all paths

$$P(x) = \sum_{\pi} P(x, \pi)$$

Prob of emissions, over all paths

Decoding

3. Viterbi decoding

$$\pi^* = \operatorname{argmax}_{\pi} P(x, \pi)$$

Most likely path

4. Posterior decoding

$$\pi^{\wedge} = \{ \pi_i \mid \pi_i = \operatorname{argmax}_k \sum_{\pi} P(\pi_i = k | x) \}$$

Path containing the most likely state at any time point.

Learning

5. Supervised learning, given π

$$\Lambda^* = \operatorname{argmax}_{\Lambda} P(x, \pi | \Lambda)$$

6. Unsupervised learning.

$$\Lambda^* = \operatorname{argmax}_{\Lambda} \max_{\pi} P(x, \pi | \Lambda)$$

Viterbi training, best path

6. Unsupervised learning

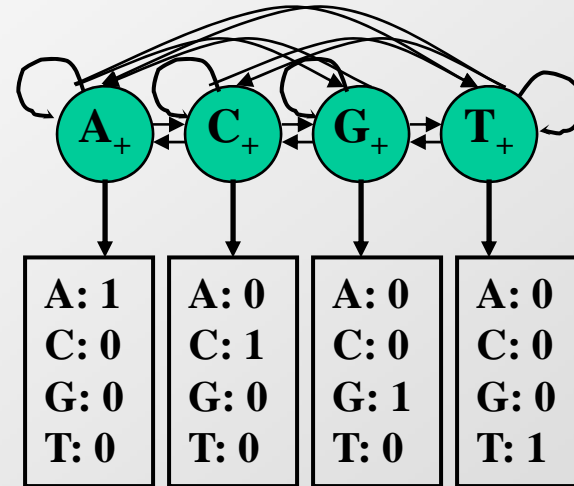
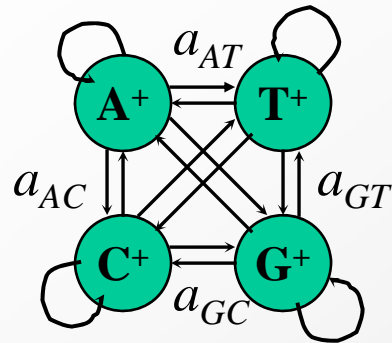
$$\Lambda^* = \operatorname{argmax}_{\Lambda} \sum_{\pi} P(x, \pi | \Lambda)$$

Baum-Welch training, over all paths

Introducing memory

- State, emissions, only depend on **current** state
- How do you count **di-nucleotide** frequencies?
 - CpG islands
 - Codon triplets
 - Di-codon frequencies
- Introducing memory to the system
 - Expanding the number of states

Example 2: CpG islands: incorporating memory



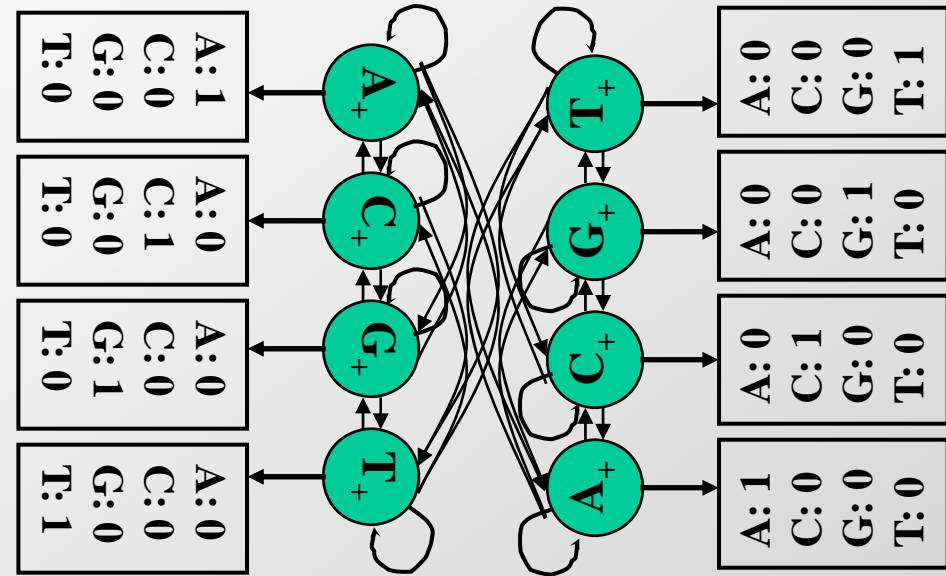
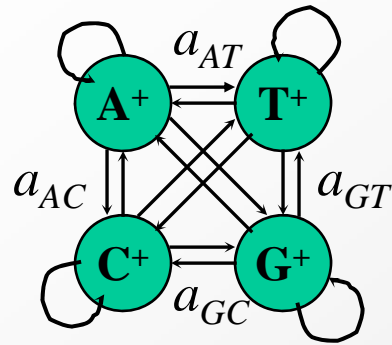
- Markov Chain

- Q: states
- p: initial state probabilities
- A: transition probabilities

- HMM

- Q: states
- V: observations
- p: initial state probabilities
- A: transition probabilities
- E: emission probabilities

Counting nucleotide transitions: Markov/HMM



- Markov Chain

- Q: states
- p: initial state probabilities
- A: transition probabilities

- HMM

- Q: states
- V: observations
- p: initial state probabilities
- A: transition probabilities
- E: emission probabilities

What have we learned ?

- Modeling sequential data
 - Recognize a **type** of sequence, genomic, oral, verbal, visual, etc...
- Definitions
 - Markov Chains
 - Hidden Markov Models (HMMs)
- Simple examples
 - Recognizing GC-rich regions.
 - Recognizing CpG dinucleotides
- Our first computations
 - Running the model: know model \rightarrow generate sequence of a 'type'
 - Evaluation: know model, emissions, states \rightarrow p?
 - Viterbi: know model, emissions \rightarrow find optimal path
 - Forward: know model, emissions \rightarrow total p over all paths
- Next time:
 - Posterior decoding
 - Supervised learning
 - Unsupervised learning: Baum-Welch, Viterbi training