

Markov Models

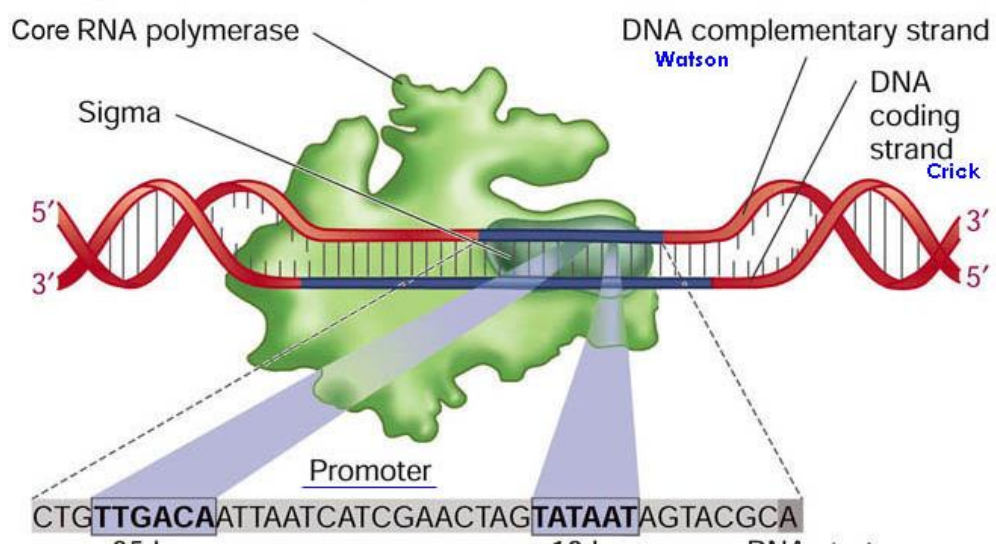
BMI/CS 576

www.biostat.wisc.edu/bmi576/

Mark Craven

craven@biostat.wisc.edu

Motivation for sequence modeling



these sequences are E. coli promoters

```
tctgaaatgagctgttgacaattaatcatcgaactagttaactagtacgcaagttca
accggaagaaaaccgtgacatTTTTAACCGTTTGTtacaaggtaaaggcgacgccgc
aaattaaaattttattgacttaggtcactaaatactttaaccaatataggcatagcg
ttgtcataatcgacttgtaaaccaaattgaaaagatttaggtttacaagtctacacc
catcctcgcaccagtcgacgacggtttacgctttacgtatagtggcgacaatttttt
tccagtataatttggttggcataattaagtacgacgagtaaaattacatacctgcccg
acagttatccactattcctgtggataaccatgtgtattagagtttagaaaacacgagg
```

these sequences are not promoters

```
atagtctcagagtcttgacctactacgccagcattttggcgggtgtaagctaaccatt
aactcaagggtgatacggcgagacttgcgagccttgtccttgcggtacacagcagcg
ttactgtgaacattattcgtctccgcgactacgatgagatgcctgagtgcttccggt
tattctcaacaagattaaccgacagattcaatctcgtggatggacgttcaacattga
aacgagtcaatcagaccgctttgactctggtattactgtgaacattattcgtctccg
aagtgcttagcttcaaggtcacggatacaccgaagcagcctcgtcctcaatggcc
gaagaccacgcctcgccaccgagtagacccttagagagcagatgtcagcctcgacaact
```

How can we tell the difference? Is this sequence a promoter?

```
ccatcaaaaaaatatttctcaacataaaaaactttgtgtaacttgtaacgctacat
```

Motivation for Markov models in computational biology

- there are many cases in which we would like to represent the statistical regularities of some class of sequences
 - genes
 - various regulatory sites in DNA (e.g. promoters)
 - proteins in a given family
 - etc.
- Markov models are well suited to this type of task

Markov chain models

- a Markov chain model is defined by
 - a set of states
 - some states *emit* symbols
 - other states (e.g. the *begin* and *end*) are *silent*
 - a set of transitions with associated probabilities
 - the transitions emanating from a given state define a distribution over the possible next states

Markov chain models

- Let X be a sequence of random variables $X_1 \dots X_L$ representing a biological sequence
- from the chain rule of probability

$$\begin{aligned} P(X) &= P(X_L, X_{L-1}, \dots, X_1) = \\ &= P(X_L | X_{L-1}, \dots, X_1) \times \\ &\quad P(X_{L-1} | X_{L-2}, \dots, X_1) \times \\ &\quad \vdots \\ &\quad P(X_1) \end{aligned}$$

Markov chain models

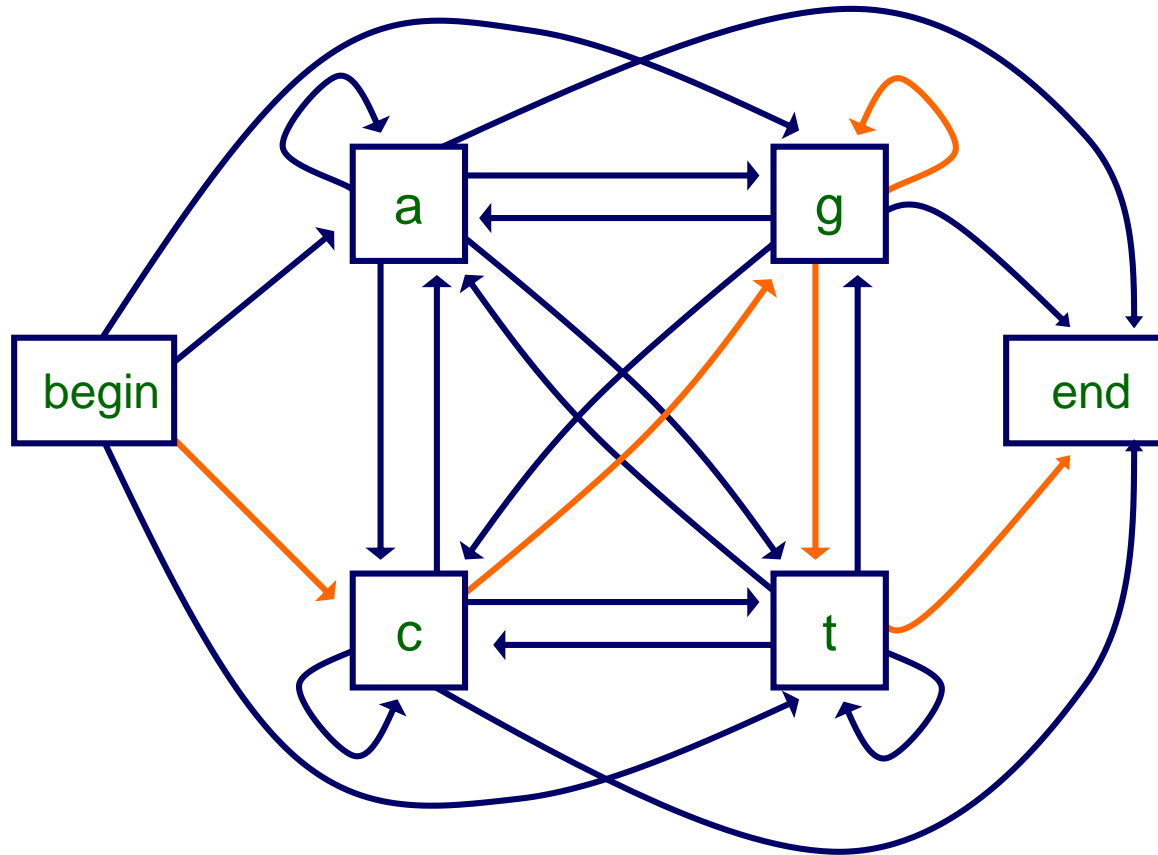
- from the chain rule we have

$$P(X) = P(X_L | X_{L-1}, \dots, X_1) P(X_{L-1} | X_{L-2}, \dots, X_1) \dots P(X_1)$$

- key property of a (1st order) Markov chain: the probability of each X_i depends only on the value of X_{i-1}

$$\begin{aligned} P(X) &= P(X_L | X_{L-1}) P(X_{L-1} | X_{L-2}) \dots P(X_2 | X_1) P(X_1) \\ &= P(X_1) \prod_{i=2}^L P(X_i | X_{i-1}) \end{aligned}$$

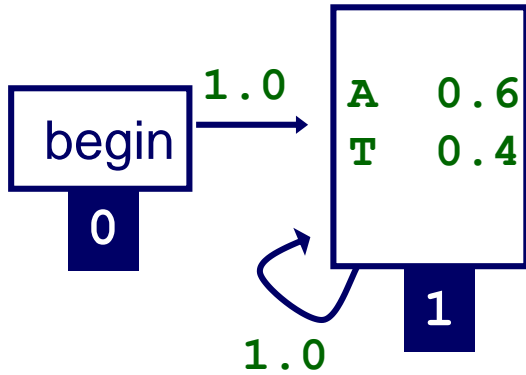
Markov chain model



The probability of a sequence $cggt$ for a given model:

$$P(cggt) = P(c)P(g|c)P(g|g)P(t|g)P(\text{end}|t)$$

Why we need an end state to define a distribution over varying length sequences



$$P(\text{A}) = 0.6$$

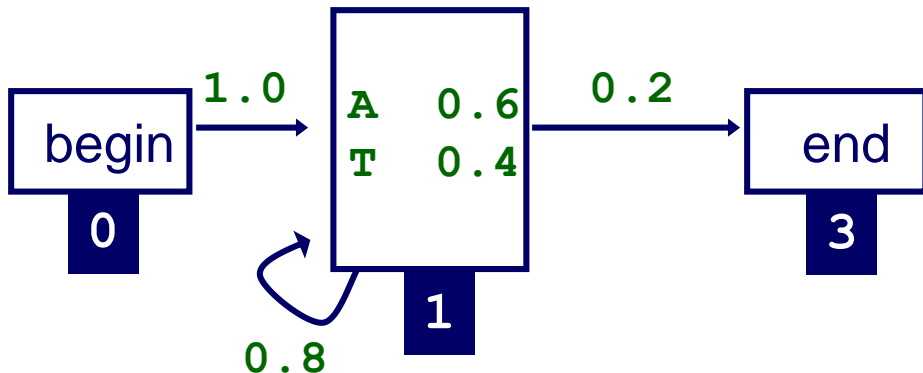
$$P(\text{T}) = 0.4$$

$$P(\text{AA}) = 0.36$$

$$P(\text{AT}) = 0.24$$

$$P(\text{TA}) = 0.24$$

$$P(\text{TT}) = 0.16$$



$$P(\text{A}) = 0.12$$

$$P(\text{T}) = 0.08$$

$$P(\text{AA}) = 0.0576$$

$$P(\text{AT}) = 0.0384$$

$$P(\text{TA}) = 0.0384$$

$$P(\text{TT}) = 0.0256$$

$$P(L=1) = 0.2$$

$$P(L=2) = 0.16$$

Estimating the model parameters

- Given some sequences, how can we determine the probability parameters of our model?
 - maximum likelihood estimation
 - Bayesian approach – regularization, priors
- estimate 1st order parameters using Laplace estimates with the sequences

gccgcgcttg

gcttggtggc

tggccgttgc

$$P(a | g) = \frac{0 + 1}{12 + 4}$$

$$P(c | g) = \frac{7 + 1}{12 + 4}$$

$$P(g | g) = \frac{3 + 1}{12 + 4}$$

$$P(t | g) = \frac{2 + 1}{12 + 4}$$

A Bayesian approach

- instead of estimating parameters strictly from the data, we could start with some prior belief for each
- in general, use a prior Dirichlet distribution as a conjugate prior to the observed multinomial data distribution
- the outcome reduces to *m*-estimates

$$P(a) = \frac{n_a + p_a m}{\left(\sum_i n_i \right) + m}$$

observed frequency of a

prior probability of a

number of “virtual” instances

iterate over all symbols/transitions

- their most simple form = *Laplace estimates*
(uniform p_a , $m=1/p_a$)

Higher order Markov chains

- we can build more “memory” into our states by using a higher order Markov model
- additional history can have predictive value
- example:
 - predict the next word in this sentence fragment
“... the__” (duck, end, grain, tide, wall, ...?)
 - now predict it given more history
“... against the __” (duck, end, grain, tide, wall, ...?)
“swim against the __” (duck, end, grain, tide, wall, ...?)

Higher order Markov chains

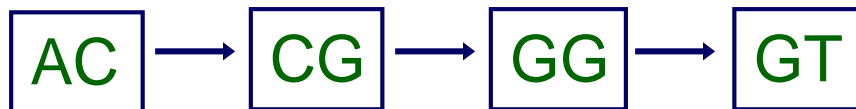
- an n th order Markov chain over some alphabet A is equivalent to a first order Markov chain over the alphabet A^n of n -tuples

- example: a 2nd order Markov model for DNA can be treated as a 1st order Markov model over alphabet

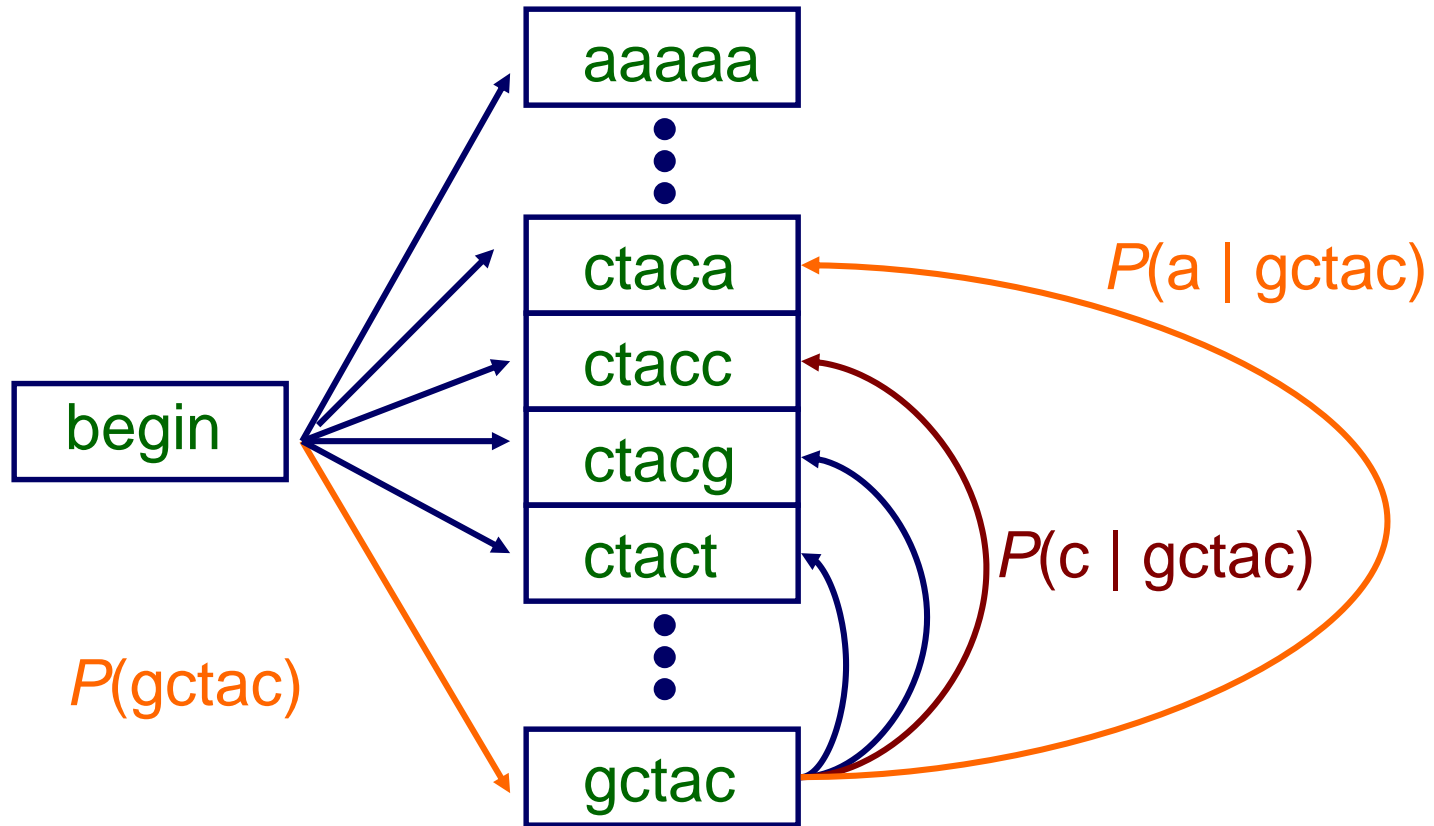
AA, AC, AG, AT, CA, CC, CG, CT, GA, GC, GG, GT,
TA, TC, TG, TT

- caveat: we process a sequence one character at a time

A C G G T



A fifth-order Markov chain

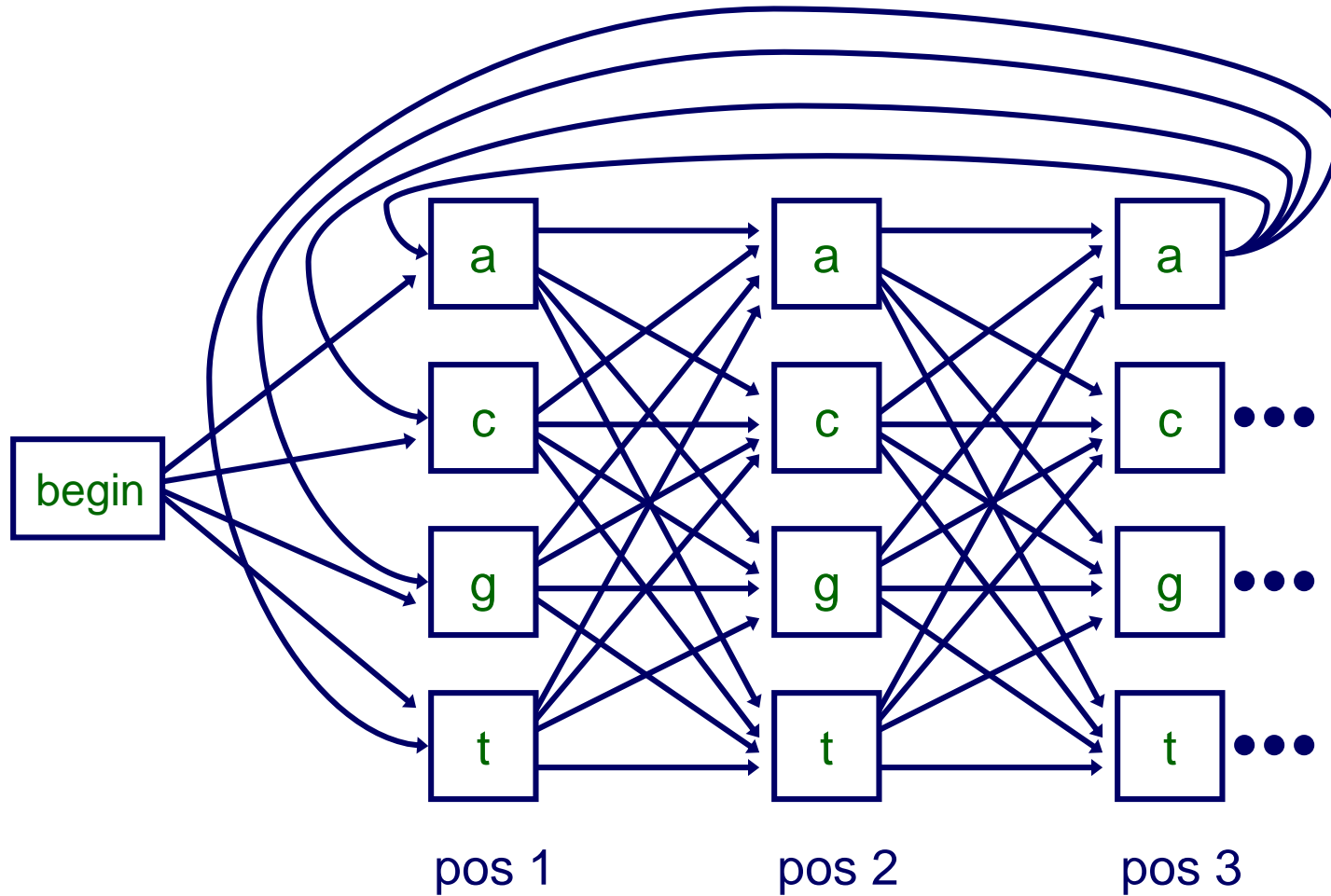


$$P(gctaca) = P(gctac)P(a | gctac)$$

Inhomogenous Markov chains

- in the Markov chain models we have considered so far, the probabilities do not depend on our position in a given sequence
- in an *inhomogeneous* Markov model, we can have different distributions at different positions in the sequence
- consider modeling codons in protein coding regions

An inhomogeneous Markov chain

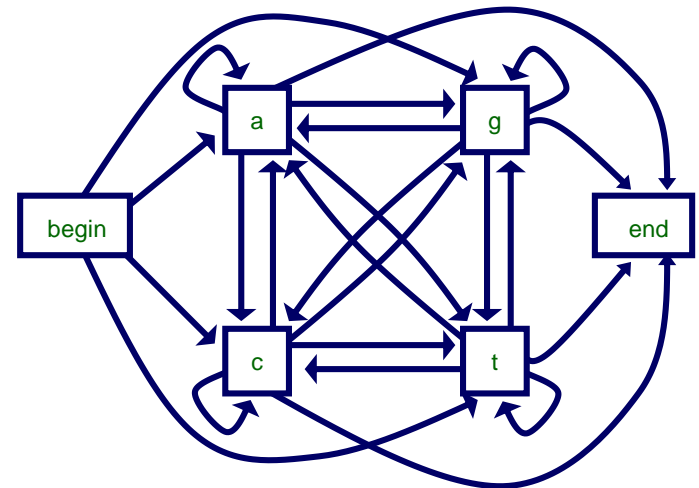
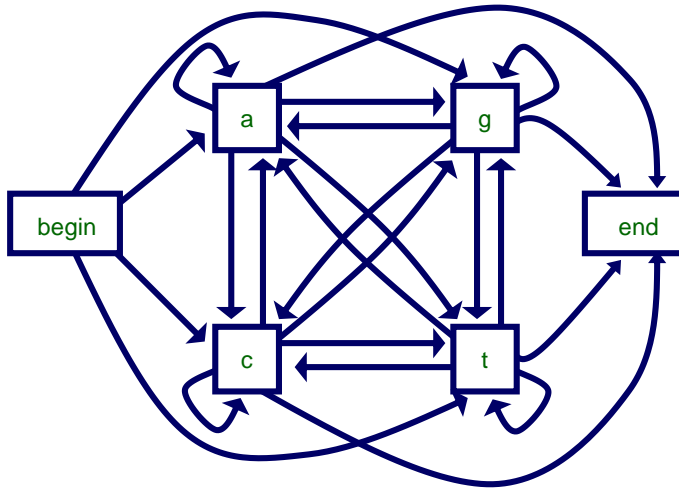


Example application

- CpG islands
 - CG dinucleotides are rarer in eukaryotic genomes than expected given the marginal probabilities of C and G
 - but the regions upstream of genes are richer in CG dinucleotides than elsewhere – *CpG islands*
 - useful evidence for finding genes
- could predict CpG islands with Markov chains
 - one to represent CpG islands
 - one to represent the rest of the genome

CpG islands as a classification task

1. train two Markov models: one to represent CpG island sequence regions, another to represent other sequence regions (*null*)



2. given a test sequence, use two models to
 - determine probability that sequence is a CpG island
 - classify the sequence (*CpG* or *null*)

Markov chains for discrimination

- parameters estimated for CpG and null models
 - human sequences containing 48 CpG islands
 - 60,000 nucleotides

$P(c | a)$

	<i>a</i>	<i>c</i>	<i>g</i>	<i>t</i>
<i>a</i>	.18	.27	.43	.12
<i>c</i>	.17	.37	.27	.19
<i>g</i>	.16	.34	.38	.12
<i>t</i>	.08	.36	.38	.18

CpG

	<i>a</i>	<i>c</i>	<i>g</i>	<i>t</i>
<i>a</i>	.30	.21	.28	.21
<i>c</i>	.32	.30	.08	.30
<i>g</i>	.25	.24	.30	.21
<i>t</i>	.18	.24	.29	.29

null

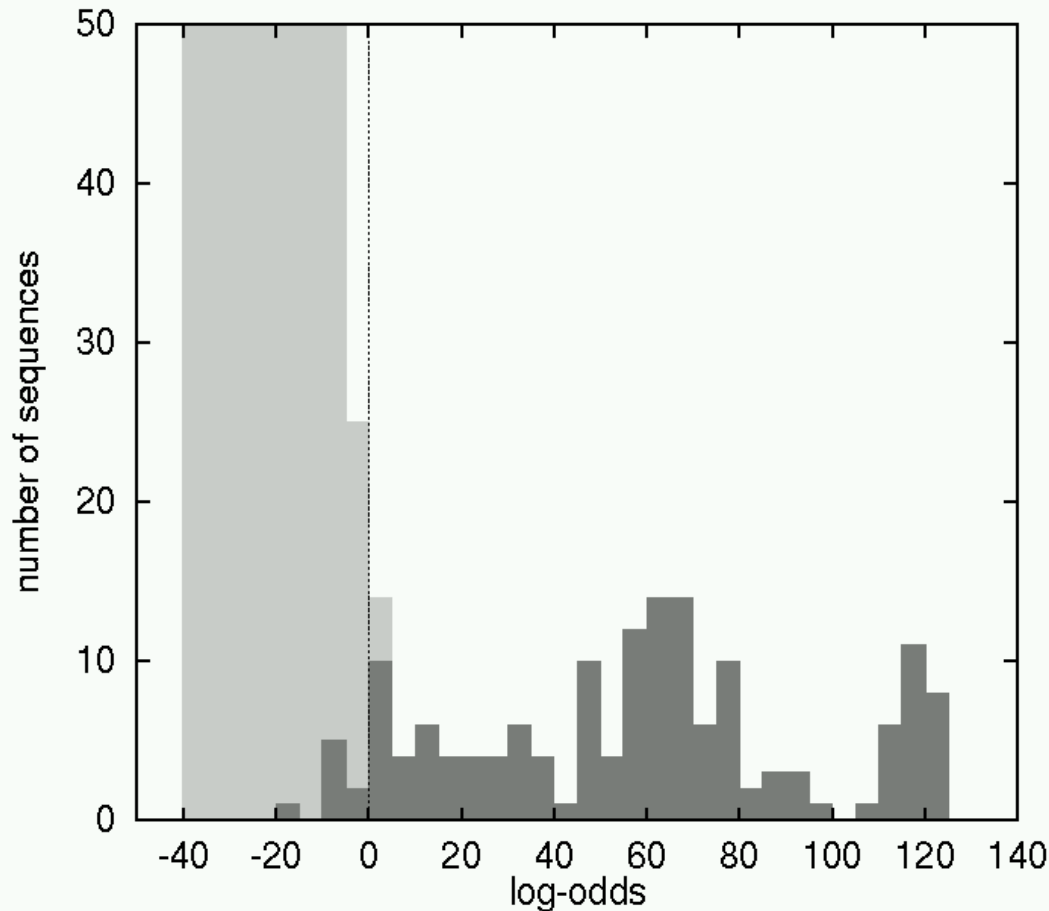
Markov chains for discrimination

- using Bayes' rule tells us

$$\begin{aligned} P(CpG | x) &= \frac{P(x | CpG)P(CpG)}{P(x)} \\ &= \frac{P(x | CpG)P(CpG)}{P(x | CpG)P(CpG) + P(x | null)P(null)} \end{aligned}$$

- if we don't take into account prior probabilities of two classes ($P(CpG)$ and $P(null)$) then we just need to compare $P(x | CpG)$ and $P(x | null)$

Markov chains for discrimination



- light bars represent negative sequences
- dark bars represent positive sequences (i.e. CpG islands)
- the actual figure here is not from a CpG island discrimination task, however

Figure from A. Krogh, "An Introduction to Hidden Markov Models for Biological Sequences" in Computational Methods in Molecular Biology, Salzberg et al. editors, 1998.

The hidden part of the problem

- in the Markov models we've considered previously, it is clear which state accounts for each part of the observed sequence
- we'll distinguish between the *observed* parts of a problem and the *hidden* parts
- in **hidden markov models**, there are multiple states that could account for each part of the observed sequence – this is the hidden part of the problem

The parameters of an HMM

- as in Markov chain models, we have transition probabilities

$$a_{kl} = P(\pi_i = l \mid \pi_{i-1} = k)$$

probability of a transition from state k to l

π represents a path (sequence of states) through the model

- since we've decoupled states and characters, we might also have emission probabilities

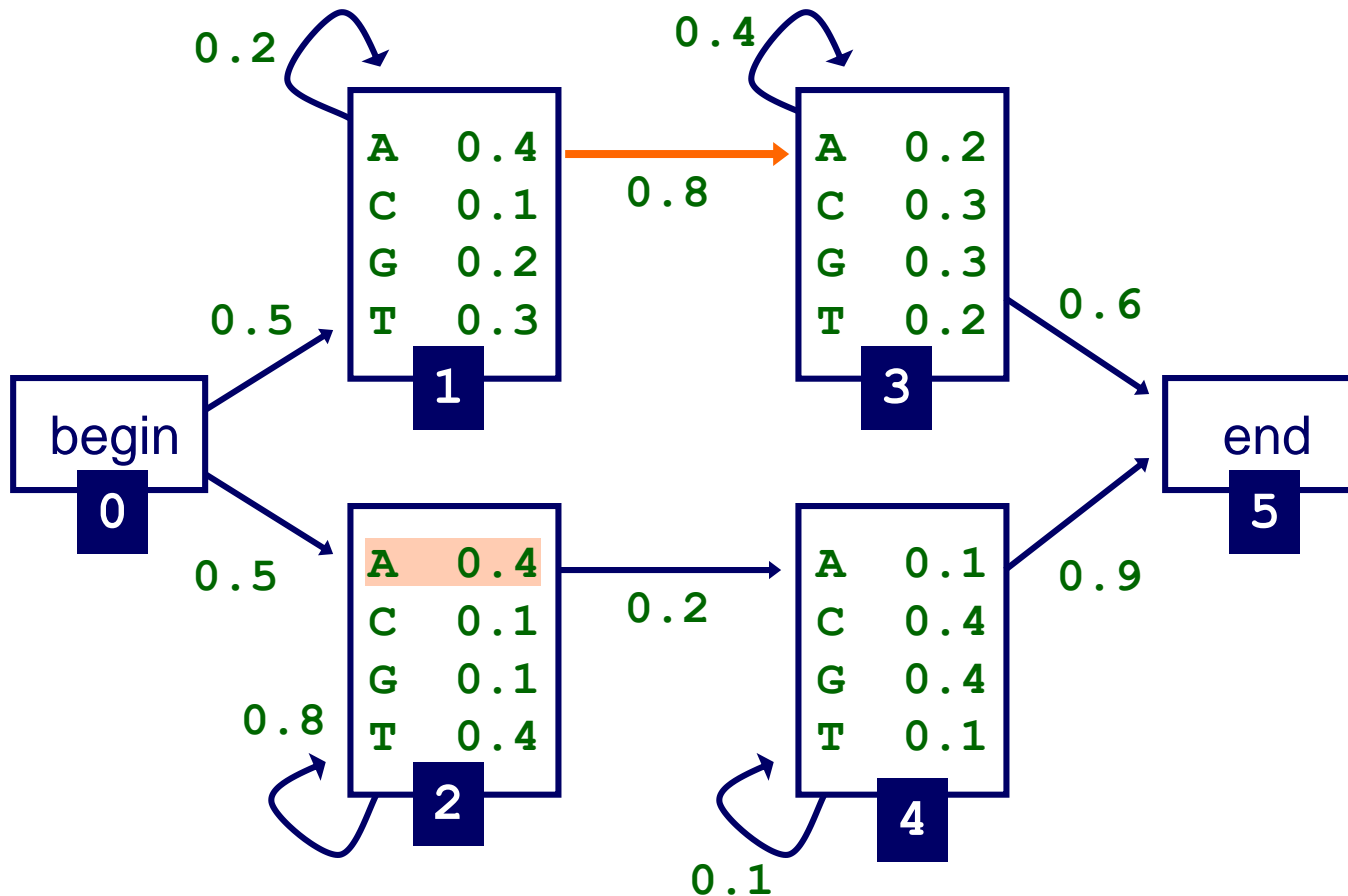
$$e_k(b) = P(x_i = b \mid \pi_i = k)$$

probability of emitting character b in state k

A simple HMM with emission parameters

a_{13} probability of a transition from state 1 to state 3

$e_2(A)$ probability of emitting character A in state 2



Three important HMM questions

- How likely is a given sequence given the model?
the Forward algorithm
- What is the most probable “path” for generating a given sequence?
the Viterbi algorithm
- How can we learn the HMM parameters given a set of sequences?
the Forward-Backward (Baum-Welch) algorithm

Learning and prediction tasks

- *learning*
 - Given:** a model, a set of training sequences
 - Do:** find model parameters that explain the training sequences with relatively high probability (goal is to find a model that *generalizes* well to sequences we haven't seen before)
- *classification*
 - Given:** a set of models representing different sequence classes, a test sequence
 - Do:** determine which model/class best explains the sequence
- *segmentation*
 - Given:** a model representing different sequence classes, a test sequence
 - Do:** segment the sequence into subsequences, predicting the class of each subsequence