

Modeling Dinucleotides

- **Di**nucleotides are important here. We need to model them in a sequence.
- Build a model that generates sequences in which the probability of a symbol depends on the previous symbol (why?).
- Markov Chain!

Markov Chain A Markov Chain is defined as:

- A set of states
- For each pair of states *i* and *j*, a transition probability a_{ij}.

 $-\Sigma_j a_{ij} = 1$

Markov Chain

- We transition from one state to another in discrete time steps *n* = 1, 2, 3, ...
- If we are at state *i* in time step *n*, we go to state *j* in time step *n*+1 with probability a_{ij}.
- The state at time *n*, *x_n*, depends on the states *x_{n-1}*, *x_{n-2}*, ... only through the most recent state *x_{n-1}*.

 $p(x_n = j \mid x_0, x_1, ..., x_{n-2}, x_{n-1} = i)$ $p(x_n = j \mid x_{n-1} = i)$



=

= a_{ii}



 $p(x_n = j \mid x_0 = i)$

The probability of being in state j at time n given that state i is the starting state:

 $p(x_n = j \mid x_0 = i) = \mathsf{P}^n_{ij}$



Probability of a path (or a sequence)

The probability of a given sequence of states $x_1...x_n$ is: $p(x_1...x_n) =$ $p(x_1...x_{n-1}, x_n) =$ $p(x_n, x_1...x_{n-1}) =$ $p(x_n | x_{n-1})p(x_1...x_{n-1}) =$ $p(x_n | x_{n-1})p(x_1...x_{n-1}) =$ $p(x_n | x_{n-1})p(x_1...x_{n-1}) =$ $p(x_1 | x_{n-1}) =$ $p(x_1) \prod_{i=2...n} a_{x^{i+1}x}$ Modeling the beginning and end of the sequence

- $p(x_1...x_n) = p(x_1)\prod_{i=2...n} a_{x^{i,i}x^i}$
- what is $p(x_1)$? Depends on how we start.
- Add a distinct start state $x_0 = S$.











Questions

• Done with Q1.

- Q2: Given a long sequence, does it contain a CpG island or not?
- How can we answer Q2?

Markov Chain model

- The Markov Chain model that we have just build can be used
- Calculate the log-odds score for windows of size, say 100, in the sequence
- CpG islands will stand out with positive values

Problems?

The previous approach is unsatisfactory

- CpG islands have variable length
- Why use a window of size 100? Why not 10 or 50 or 200? (no way to tell best size, could be average length of CpG island, but still unsatisfactory)



Better solution

- Represent CpG islands and non CpG islands in one model
- Both Markov chains build earlier put together with small transition probability between them
- We will have two states for each nucleotides → rename them A+, C+, G+, T+ and A-, C-, G-, T-



What is the big difference now?

- There is not a one-to-one correspondence between the states and the symbols.
 Given a symbol C, it could have been generated by state C+ or state C-
- Before, a sequence uniquely determines the path
- Now, for a given sequence, we want to find the most likely path



Hidden Markov Model

A Hidden Markov model HMM is defined as: [state is hidden, decouple states from symbols]

- A set of hidden states
- For each pair of states i and j, a transition probability a_{ij}.
- $-\Sigma_i a_{ii} = 1$
- For each state k, emission probabilities $e_k(b) = p(x_i = b \mid \pi_i = k)$ [now we use variable π for states and variable x for symbols]
- $-\Sigma_b e_k(b) = 1$ for each state k
- Markov property: $p(\pi_n = j \mid x_0...x_{n-1}, \pi_0...\pi_{n-2}, \pi_{n-1} = i) = p(\pi_n = j \mid \pi_{n-1} = i) = a_{ij}$

Questions with HMMs

- *Evaluation*: given *x*, what is the probability *p*(*x*) that it was produced by the model?
- *Decoding*: given *x*, what is the most probable path that produces *x* in the model?
- Learning: given x, what are the most probable parameters (transitional probabilities and emission probabilities) of the model?

The dishonest Casino

- Casino uses a fair die most of the time
- Loaded die has p(1) = p(2) = p(3) = p(4) = p(5) = 0.1, p(6) = 0.5
- Casino switches from fair to loaded with a probability of $0.05\,$
- Switches back with probability 0.1

[think about similarities with CpG island]



Most probable path
Label start state and end state by 0.
 The joint probability of observing a sequence of symbols x = x₁x_n emitted by a sequence of states π = π₁π_n; p(x,π) = a_{0π1}.e_{π1}(x₁)a_{πn-1} n_ne_{πn}(x_n).a_{πn0} = a_{0π1}Π_{μ=1n} e_{πn}(x₁)a_{πi}π_{i+1} where π_{n+1} = 0 We want to find π* = argmax_π p(π x) = argmax_π p(x,π) Try all possible π: EXPONENTIAL!
Sad Memory







