

Students:

- Julian Buetzer
- Giacomo Frison
- Cindy Guerrero

Viterbi Algorithm: Manual Solution for CG-rich Genome Regions

Problem Setup: Case 1

Sequence: GGACTGAA (8 nucleotides)

Hidden States: - **R** (CG-rich): emits C and G with higher probability - **P** (CG-poor): emits T and A with higher probability

Initial Probabilities: - $_R = 0.5$ - $_P = 0.5$

Transition Probabilities: - $a_{RR} = 0.5$ ($R \rightarrow R$) - $a_{RP} = 0.5$ ($R \rightarrow P$) - $a_{PR} = 0.6$ ($P \rightarrow R$) - $a_{PP} = 0.4$ ($P \rightarrow P$)

Emission Probabilities:

Nucleotide	State R	State P
T	0.2	0.3
C	0.3	0.2
A	0.2	0.3
G	0.3	0.2

Viterbi Algorithm Steps

(Exercise 1) Use a Viterbi algorithm to define the most likely sequence of hidden states for sequence GGACTGAA.

The Viterbi algorithm finds the most likely sequence of hidden states using dynamic programming.

Notation

- $P[t][s]$: Maximum probability of ending at state s at time t
- $Q[t][s]$: Previous state that led to state s at time t (for backtracking)

Step 0: Initialization ($t=0$, observation: G)

$$P[0][R] = _R \times b_R(G) = 0.5 \times 0.3 = \mathbf{0.15}$$

$$P[0][P] = _P \times b_P(G) = 0.5 \times 0.2 = \mathbf{0.10}$$

$Q[0][R] = \text{None}$ (initial state)
 $Q[0][P] = \text{None}$ (initial state)

Step 1: t=1 (observation: G)

For state R: - From R: $P[0][R] \times a_{RR} \times b_R(G) = 0.15 \times 0.5 \times 0.3 = 0.0225$ - From P: $P[0][P] \times a_{PR} \times b_R(G) = 0.10 \times 0.6 \times 0.3 = 0.0180$ - $P[1][R] = \max(0.0225, 0.0180) = \mathbf{0.0225}$ - $Q[1][R] = R$

For state P: - From R: $P[0][R] \times a_{RP} \times b_P(G) = 0.15 \times 0.5 \times 0.2 = 0.0150$ - From P: $P[0][P] \times a_{PP} \times b_P(G) = 0.10 \times 0.4 \times 0.2 = 0.0080$ - $P[1][P] = \max(0.0150, 0.0080) = \mathbf{0.0150}$ - $Q[1][P] = R$

Step 2: t=2 (observation: A)

For state R: - From R: $P[1][R] \times a_{RR} \times b_R(A) = 0.0225 \times 0.5 \times 0.2 = 0.00225$ - From P: $P[1][P] \times a_{PR} \times b_R(A) = 0.0150 \times 0.6 \times 0.2 = 0.00180$ - $P[2][R] = \max(0.00225, 0.00180) = \mathbf{0.00225}$ - $Q[2][R] = R$

For state P: - From R: $P[1][R] \times a_{RP} \times b_P(A) = 0.0225 \times 0.5 \times 0.3 = 0.003375$ - From P: $P[1][P] \times a_{PP} \times b_P(A) = 0.0150 \times 0.4 \times 0.3 = 0.00180$ - $P[2][P] = \max(0.003375, 0.00180) = \mathbf{0.003375}$ - $Q[2][P] = R$

Step 3: t=3 (observation: C)

For state R: - From R: $P[2][R] \times a_{RR} \times b_R(C) = 0.00225 \times 0.5 \times 0.3 = 0.0006075$ - From P: $P[2][P] \times a_{PR} \times b_R(C) = 0.003375 \times 0.6 \times 0.3 = 0.0006075$ - $P[3][R] = \max(0.0006075, 0.0006075) = \mathbf{0.0006075}$ - $Q[3][R] = P$

For state P: - From R: $P[2][R] \times a_{RP} \times b_P(C) = 0.00225 \times 0.5 \times 0.2 = 0.000225$ - From P: $P[2][P] \times a_{PP} \times b_P(C) = 0.003375 \times 0.4 \times 0.2 = 0.00027$ - $P[3][P] = \max(0.000225, 0.00027) = \mathbf{0.00027}$ - $Q[3][P] = P$

Step 4: t=4 (observation: T)

For state R: - From R: $P[3][R] \times a_{RR} \times b_R(T) = 0.0006075 \times 0.5 \times 0.2 = 0.0006075$ - From P: $P[3][P] \times a_{PR} \times b_R(T) = 0.00027 \times 0.6 \times 0.2 = 0.000324$ - $P[4][R] = \max(0.0006075, 0.000324) = \mathbf{0.0006075}$ - $Q[4][R] = R$

For state P: - From R: $P[3][R] \times a_{RP} \times b_P(T) = 0.0006075 \times 0.5 \times 0.3 = 0.00091125$ - From P: $P[3][P] \times a_{PP} \times b_P(T) = 0.00027 \times 0.4 \times 0.3 =$

$$0.0000324 - \mathbf{P[4][P]} = \max(0.000091125, 0.0000324) = \mathbf{0.000091125} - \mathbf{Q[4][P]} \\ = R$$

Step 5: t=5 (observation: G)

For state R: - From R: $P[4][R] \times a_{RR} \times b_R(G) = 0.00006075 \times 0.5 \times 0.3 = 0.0000091125$ - From P: $P[4][P] \times a_{PR} \times b_R(G) = 0.000091125 \times 0.6 \times 0.3 = 0.0000164025$ - $\mathbf{P[5][R]} = \max(0.0000091125, 0.0000164025) = \mathbf{0.0000164025}$ - $\mathbf{Q[5][R]} = P$

For state P: - From R: $P[4][R] \times a_{RP} \times b_P(G) = 0.00006075 \times 0.5 \times 0.2 = 0.000006075$ - From P: $P[4][P] \times a_{PP} \times b_P(G) = 0.000091125 \times 0.4 \times 0.2 = 0.00000729$ - $\mathbf{P[5][P]} = \max(0.000006075, 0.00000729) = \mathbf{0.00000729}$ - $\mathbf{Q[5][P]} = P$

Step 6: t=6 (observation: A)

For state R: - From R: $P[5][R] \times a_{RR} \times b_R(A) = 0.0000164025 \times 0.5 \times 0.2 = 0.00000164025$ - From P: $P[5][P] \times a_{PR} \times b_R(A) = 0.00000729 \times 0.6 \times 0.2 = 0.0000008748$ - $\mathbf{P[6][R]} = \max(0.00000164025, 0.0000008748) = \mathbf{0.00000164025}$ - $\mathbf{Q[6][R]} = R$

For state P: - From R: $P[5][R] \times a_{RP} \times b_P(A) = 0.0000164025 \times 0.5 \times 0.3 = 0.000002460375$ - From P: $P[5][P] \times a_{PP} \times b_P(A) = 0.00000729 \times 0.4 \times 0.3 = 0.0000008748$ - $\mathbf{P[6][P]} = \max(0.000002460375, 0.0000008748) = \mathbf{0.000002460375}$ - $\mathbf{Q[6][P]} = R$

Step 7: t=7 (observation: A)

For state R: - From R: $P[6][R] \times a_{RR} \times b_R(A) = 0.00000164025 \times 0.5 \times 0.2 = 0.000000164025$ - From P: $P[6][P] \times a_{PR} \times b_R(A) = 0.000002460375 \times 0.6 \times 0.2 = 0.000000295245$ - $\mathbf{P[7][R]} = \max(0.000000164025, 0.000000295245) = \mathbf{0.000000295245}$ - $\mathbf{Q[7][R]} = P$

For state P: - From R: $P[6][R] \times a_{RP} \times b_P(A) = 0.00000164025 \times 0.5 \times 0.3 = 0.0000002460375$ - From P: $P[6][P] \times a_{PP} \times b_P(A) = 0.000002460375 \times 0.4 \times 0.3 = 0.000000295245$ - $\mathbf{P[7][P]} = \max(0.0000002460375, 0.000000295245) = \mathbf{0.000000295245}$ - $\mathbf{Q[7][P]} = P$

Backtracking to Find the Most Likely Path

At the final time step ($t=7$): - $P[7][R] = 0.000000295245$ - $P[7][P] = 0.000000295245$

Both states have the same probability! We can choose either, but typically we choose the one with the higher index or use a tie-breaking rule. At this point, we can trace back from P (since it's the last state in our Q matrix, or A matrix according to the other book, Lecture 5, part II, slide 24):

Path reconstruction: - $t=7: P$ (final state) - $t=6: Q[7][P] = P$ - $t=5: Q[6][P] = R$ - $t=4: Q[5][R] = P$ - $t=3: Q[4][P] = R$ - $t=2: Q[3][R] = P$ - $t=1: Q[2][P] = R$ - $t=0: Q[1][R] = R$

Most likely sequence of hidden states: R-R-P-R-P-R-P-P

Sequence alignment:

Observation: G G A C T G A A
Hidden State: R R P R P R P P

Summary

The Viterbi algorithm has determined that the most likely sequence of hidden states for "GGACTGAA" is **R-R-P-R-P-R-P-P**, where: - **R** = CG-rich region - **P** = CG-poor region

This makes biological sense because: - The sequence starts with GG (two Gs), which favors the CG-rich state - The middle contains A, C, T which are more likely in CG-poor regions - The algorithm balances emission probabilities with transition probabilities to find the optimal path

Probability $P(X)$ of the sequence GGCA calculated through forward and backward algorithm:

For states, initial and transition probabilities and emission probabilities see above

Forward algorithm for sequence GGCA:

initialization ($t=1, G$) $a1(R) = R \times P(G|R) = 0.5 \times 0.3 = 0.15$ $a1(P) = P \times P(G|P) = 0.5 \times 0.2 = 0.10$

recursion ($t=2, G$) $a2(R) = [a1(R)(a_PR) + a1(P)(a_PP)] \times P(G|R) = [(0.15 \times 0.6) + (0.10 \times 0.4)] \times 0.3 = \mathbf{0.0039}$ $a2(R) = [a1(R)(a_PP) + a1(P)(a_PR)] \times P(G|R) = [(0.15 \times 0.4) + (0.10 \times 0.6)] \times 0.2 = \mathbf{0.024}$

recursion ($t=3 C$) $a3(R) = [a2(R)(a_PR) + a2(P)(a_PP)] \times P(C|R) = [(0.0039 \times 0.6) + (0.024 \times 0.4)] \times 0.3 = \mathbf{0.0099}$ $a3(P) = [a2(R)(a_PP) + a2(P)(a_PR)] \times P(C|P) = [(0.0039 \times 0.4) + (0.024 \times 0.6)] \times 0.2 = \mathbf{0.006}$

recursion (t=4, A) $a4(R) = [a3(R)(a_PR) + a3(P)(a_PP)] \times P(A|R) = [(0.0099 \times 0.6) + (0.006 \times 0.4)] \times 0.2 = \mathbf{0.001668}$ $a4(P) = [a3(R)(a_PR) + a3(P)(a_PP)] \times P(A|P) = [(0.0099 \times 0.4) + (0.006 \times 0.6)] \times 0.2 = \mathbf{0.001512}$

termination P(GGCA) $a4(R) + a4(P) = 0.001668 + 0.001512 = \mathbf{0.00318}$

Backward algorithm

Initialization (t=4, A) $b4(R) = 1 ; b4(P) = 1$

Step t=3, C to A $b3(R) = [0.6 \times ((P(A|R) \times b4(R)) + (0.4 \times (P(A|P) \times b4(P))) = [(0.6 \times 0.2 \times 1) + (0.4 \times 0.2 \times 1)] = 0.12 + 0.08 = 0.20$ $b3(P) = [0.6 \times ((P(A|P) \times b4(P)) + (0.4 \times (P(A|R) \times b4(R))) = [(0.6 \times 0.2 \times 1) + (0.4 \times 0.2 \times 1)] = 0.12 + 0.08 = 0.20$

Step t=2, G to C $b2(R) = [0.6 \times ((P(C|R) \times b3(R)) + (0.4 \times (P(C|P) \times b3(P))) = [(0.6 \times 0.3 \times 0.20) + (0.4 \times 0.2 \times 0.20)] = 0.036 + 0.016 = 0.052$ $b2(P) = [0.6 \times ((P(C|P) \times b3(P)) + (0.4 \times (P(C|R) \times b3(R))) = [(0.6 \times 0.2 \times 0.20) + (0.4 \times 0.3 \times 0.20)] = 0.024 + 0.024 = 0.048$

Step t=1, G to G $b1(R) = [0.6 \times ((P(G|R) \times b3(R)) + (0.4 \times (P(G|P) \times b3(P))) = [(0.6 \times 0.3 \times 0.052) + (0.4 \times 0.2 \times 0.0048)] = 0.00936 + 0.00384 = 0.0132$ $b1(P) = [0.6 \times ((P(G|P) \times b3(P)) + (0.4 \times (P(G|R) \times b3(R))) = [(0.6 \times 0.3 \times 0.0048) + (0.4 \times 0.2 \times 0.0052)] = 0.00576 + 0.00624 = 0.012$

P sequence $P(s) \times P(G|R-P) \times b1(R-P) = (0.5 \times 0.3 \times 0.0132) + (0.5 \times 0.2 \times 0.012) = 0.00198 + 0.0012 = \mathbf{0.00318}$

Scaling variables

$s = 1 / (\text{at}(S))$

Initialization t=1, G

$a1(R) + a1(P) = 0.15 + 0.10 = 0.25$ $s1 = 1/0.25 = 4$ $as1(R) = 0.6$ $as1(P) = 0.4$

t=2, G $as2(R) = [a1(R)(a_PR) + a1(P)(a_PP)] \times P(G|R) = \mathbf{0.0039}$ $as2(P) = [a1(P)(a_PR) + a1(R)(a_PP)] \times P(G|R) = \mathbf{0.024}$ $s2 = 1/(0.0039 + 0.0024) = 1/0.0063 = 15.87$ $as2(R) = 0.0039 \times 15.87 = 0.619$ $as2(P) = 0.024 \times 15.87 = 0.381$

t=3, C $as3(R) = [a2(R)(a_PR) + a2(P)(a_PP)] \times P(G|R) = \mathbf{0.0099}$ $as3(P) = [a2(P)(a_PR) + a2(R)(a_PP)] \times P(G|R) = \mathbf{0.006}$ $s3 = 1/(0.0099 + 0.006) = 1/0.0159 = 62.9$ $as3(R) = 0.0099 \times 62.9 = 0.622$ $as3(P) = 0.006 \times 62.9 = 0.378$

t=4, A $as4(R) = [a3(R)(a_PR) + a3(P)(a_PP)] \times P(G|R) = \mathbf{0.001668}$ $as4(P) = [a3(P)(a_PR) + a3(R)(a_PP)] \times P(G|R) = \mathbf{0.001512}$ $s4 = 1/(0.001668 + 0.001512) = 1/0.00318 = 314.5$ $as4(R) = 0.001668 \times 314.5 = 0.525$ $as4(P) = 0.001512 \times 314.5 = 0.475$

$$\mathbf{P(GGCA)} \; 1 \; / \; (\text{s1} + \text{s2} + \text{s3} + \text{s4}) = 1 \; / \; (4 + 15.87 + 62.9 + 314.5) = \mathbf{0.00318}$$