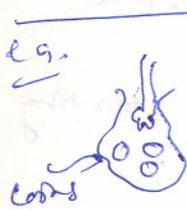
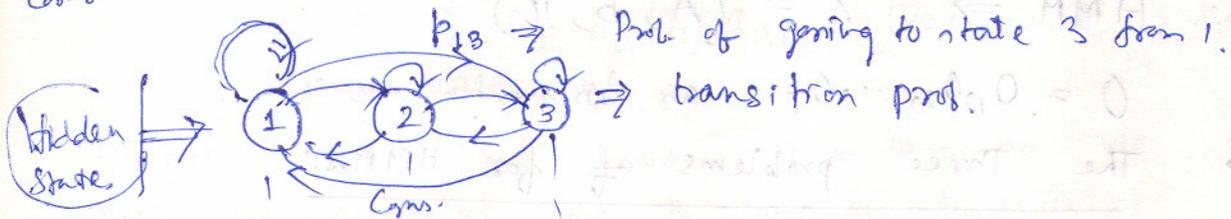


Gene Prediction:

Hidden Markov Models (HMM)



- Method.
- ① Draw a coin randomly
 - ② Toss it



$$\left. \begin{array}{ccc} p_1(H), & p_2(H), & p_3(H) \\ p_1(T), & p_2(T), & p_3(T) \end{array} \right\} \rightarrow \text{Individual bias.}$$

$$H/T \quad H/T \quad H/T \rightarrow \text{O/p symbols.}$$

O/p sequence depends very much on the individual bias (Prob. of an O/p symbol given a state), the transition probabilities (Prob. of going to state i from j) between various states, as well as on which state is chosen to begin the observation (Prob. of state i as initial state).

Notations:

$S = \{1, 2, \dots, N\} \Rightarrow$ set of N discrete states.

$V = \{v_1, v_2, \dots, v_M\} \Rightarrow$ set of M distinct observation symbols.

$i_t =$ state at time t .

$O_t =$ observation symbol at time t .

$A = \{a_{ij}\} \Rightarrow$ set of transition prob. i.e.

$$a_{ij} = P(i_{t+1} = j \mid i_t = i)$$

Assume a_{ij} 's are independent of time (or time invariant).

$$B = \{ b_j(k) \mid b_j(k) = P(v_k \text{ at } t \mid i_t = j) \},$$

Prob. of observation symbol v_k at state

$$\pi = \{ \pi_i \mid \pi_i = P(i_1 = i) \}, \text{ Prob. that the starting state is } i.$$

$$\therefore \text{HMM} \Rightarrow \lambda = (A, B, \pi)$$

$O = O_1, O_2, \dots, O_T$ as observation sequence.

- The Three problems of for HMMs:

Prob. 1. Given the model, $\lambda = (A, B, \pi)$ how do we compute $P(O | \lambda)$ i.e. the prob. of occurrence of the observation sequence

$$O = O_1, O_2, \dots, O_T.$$

Prob. 2: Given $\lambda = (A, B, \pi)$, how do we choose a state sequence $I = i_1, i_2, \dots, i_T$ so that $P(O, I | \lambda)$ is maximised.

$$P(O, I | \lambda) = \text{joint prob. of } (\langle O_1, i_1 \rangle, \langle O_2, i_2 \rangle, \dots, \langle O_T, i_T \rangle) \text{ given } \lambda.$$

Prob. 3: How do we derive an HMM model $\lambda = (A, B, \pi)$ so that $P(O | \lambda)$ (or $P(O, I | \lambda)$) is maximized.

The first two are analysis problems & the last one is a synthesis problem.

• Prob. L: $P(I|\lambda) = \pi_{i_1} a_{i_1 i_2} a_{i_2 i_3} \dots a_{i_{T-1} i_T}$

$$P(O|I, \lambda) = b_{i_1}(O_1) b_{i_2}(O_2) \dots b_{i_T}(O_T)$$

$$\therefore P(O|\lambda) = \sum_{\forall I} P(O|I, \lambda) P(I|\lambda)$$

$$= \sum_{\forall I = \{i_1, i_2, \dots, i_T\}} \pi_{i_1} b_{i_1}(O_1) a_{i_1 i_2} b_{i_2}(O_2) \dots a_{i_{T-1} i_T} b_{i_T}(O_T)$$

However, $|I| = N^T$, which is a large no. we need an efficient way to compute $P(O|\lambda)$.

Time Order: $\approx 2TN^2$

Forward Procedure:

Let, $\alpha_t(i) = P(O_1, O_2, \dots, O_t, O_t = i | \lambda)$

Compute $\alpha_t(i)$ inductively:

1. $\alpha_1(i) = \pi_i b_i(O_1)$, $1 \leq i \leq N$ $\Rightarrow N$ multiplications

2. for $t=1, 2, \dots, T-1$, $1 \leq j \leq N$,

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(O_{t+1}) \Rightarrow N+1 \text{ multiplications}$$

\Rightarrow for $\forall i$, $N(N+1)$ multiplications

$$\therefore P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$

$\forall t, j$, $(T-1)N(N+1)$

Time Complexity: $N + N(N+1)(T-1) \approx O(N^2 T)$

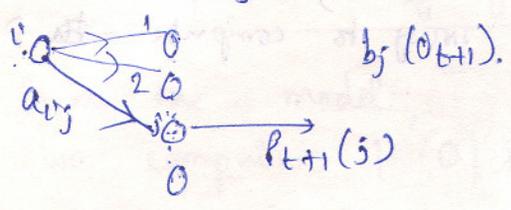
Backward Procedure:

$$\beta_t(i) = P(O_{t+1} O_{t+2} \dots O_T | \tau_t = i, \lambda)$$

1. $\beta_T(i) = 1, \quad 1 \leq i \leq N$

2. for $t = T-1, T-2, \dots, 1, \quad 1 \leq i \leq N$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

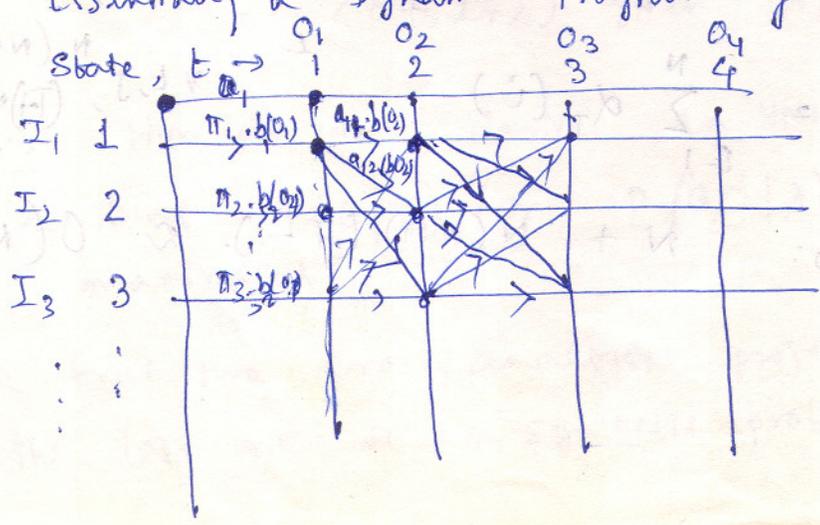


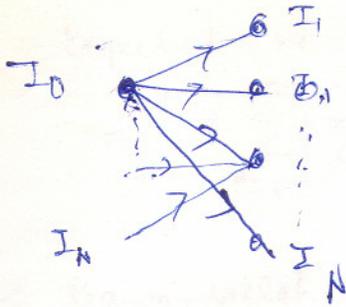
3. $P(O|\lambda) = \sum_{i=1}^N \pi_i b_i(O_1) \beta_1(i)$

© Prob. 2. Given O, λ ,
Find I to maximise $P(O, I | \lambda)$.

Viterbi Algorithm (Viterbi, 1967)

Essentially a Dynamic Programming approach.





At a node

$$s(I_i, t) = \text{maximum prob score vs ordered}$$

$$\therefore \text{at } (I_i, t+1)$$

$$= \max_{\forall j} \left\{ s(I_j, t) \cdot a_{ji} \cdot b_i(O_{t+1}) \right\}$$

1. Initial score: $s(I_i, 1) = \pi_i b_i(O_1), \quad \forall i \leq N$

2. For $t = 2, 3, \dots, T$.

$$s(I_i, t) = \max_{\forall j} \left(s(I_j, t-1) \cdot a_{ji} \cdot b_i(O_t) \right)$$

$$\text{prev. state } (\hat{i}_j, t) = \underset{j}{\text{argmax}} \left(\dots \right)$$

3. Final state:

$$P(O, I | \lambda) = \max_j \left(s(I_j, T) \right)$$

$$\text{Final-state} = \underset{j}{\text{argmax}} \left(s(I_j, T) \right)$$

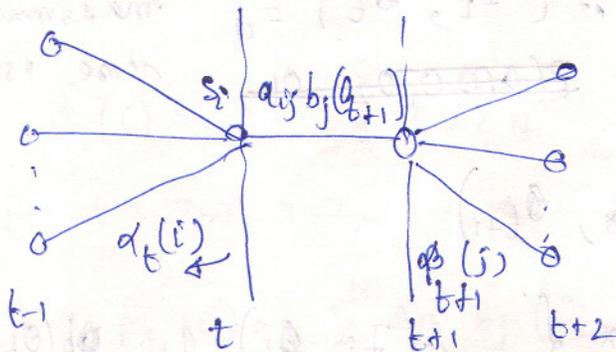
As, prob expressed in terms of products, the cost function can be conveniently expressed by logarithm to bring summation during maximization op.

$$\therefore s(I_i, 1) = \ln(\pi_i b_i(O_1))$$

$$s(I_i, t) = \max_j \left\{ s(I_j, t-1) + \ln(a_{ji} b_i(O_t)) \right\}$$

Solution to Prob. 3.

Prob. 3 Baum-Welch method.



Let, $\xi_t(i, j) = P(O_t = i, O_{t+1} = j | O, \lambda)$

$$\begin{aligned} \therefore \xi_t(i, j) &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)} \end{aligned}$$

Note, $P(O_t = i, O_{t+1} = j, O | \lambda) = \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$

Let, $\gamma_t(i) = \text{Prob.} \left(O_t = i | O, \lambda \right) = \frac{\alpha_t(i) \beta_t(i)}{P(O | \lambda)}$

$$\therefore \gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

\therefore Expected no. of transitions from i or expected no. of times that state i is visited

$$= \sum_{t=1}^{T-1} \gamma_t(i) \quad \left[\text{exclude time } t \right]$$

Expected no. of transitions from state i to j

$$= \sum_{t=1}^{T-1} \xi_t(i, j)$$

← Baum-Welch reestimation formulae:

$$\hat{\pi}_i = \gamma_{t=1}(i), \quad 1 \leq i \leq N$$

= expected no. of times state i occurs at $t=1$.

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

= (expected no. of transitions to state j for state i).

$$\hat{b}_{ij}(O_k) = \frac{\sum_{t=1}^T \gamma_{ij}^t(i) \delta_{ij}^t(O_k)}{\sum_{t=1}^T \gamma_{ij}^t(i)}$$

not. $O_{kb} = V_k$

$$\sum_{t=1}^{T-1} \gamma_t(i)$$

Alg. 1. Initialise randomly to $\lambda = (A, B, \pi)$.

2. Use reestimation till it converges.

Baum & Welch proved that the convergence is guaranteed v.p.

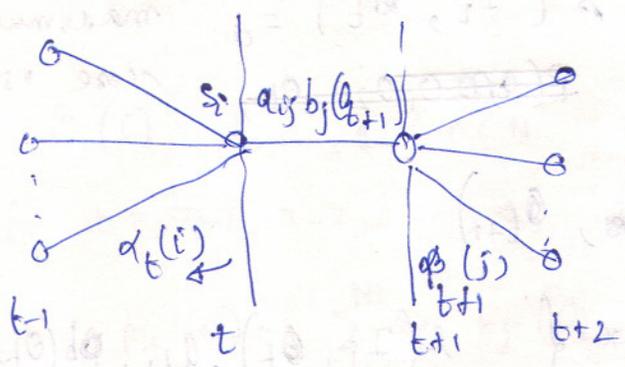
$$P(O | \lambda_{v+1}) \geq P(O | \lambda_v)$$

after v th iterat.

$v = 0, 1, \dots$

Solution to Prob. 3.

Baum-Welch method.



Let, $\xi_t(i, j) = P(x_t = i, x_{t+1} = j | O, \lambda)$

$$\begin{aligned} \therefore \xi_t(i, j) &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{P(O | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)} \end{aligned}$$

Note, $P(x_t = i, x_{t+1} = j, O | \lambda) = \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$

Let, $\gamma_t(i) = \text{Prob. } (x_t = i | O, \lambda) = \frac{\alpha_t(i) \beta_t(i)}{P(O | \lambda)}$

$$\therefore \gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

\therefore Expected no. of transitions from state i or expected no. of times that state i is visited

$$= \sum_{t=1}^{T-1} \gamma_t(i) \quad \left[\text{exclude time step } T \right]$$

Expected no. of transitions from state i to j

$$= \sum_{t=1}^{T-1} \bar{z}_t(i, j)$$

← Baum-welch reestimation formulae:

$$\hat{\pi}_i = \gamma_{t=1}(i), \quad 1 \leq i \leq N$$

= expected no. of times state i occurs at $t=1$.

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \bar{z}_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}$$

= expected no. of transitions to state j for state i .

$$\hat{b}_{ij}(O_k) = \frac{\sum_{t=1}^T \gamma_t^i(i) \delta_{ij}(O_k)}{\sum_{t=1}^T \gamma_t^i(i)}$$

not. $O_{kb} = V_k$

$$\sum_{t=1}^T \gamma_t^i(i)$$

Alg. 1. Initialize randomly to $\lambda = (A, B, \pi)$.

2. Use reestimation till at convergence.

Baum & Welch proved that the convergence is guaranteed v.p.

$$P(O | \lambda_{v+1}) > P(O | \lambda_v)$$

after v th iterat.

$v = 0, 1, \dots$

Ref. L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition", Proc. of the IEEE, vol. 77, no. 2, Feb., 1989, p. 257 - 281.

④ Explicit State duration Density in HMMs:

With conventional HMM:



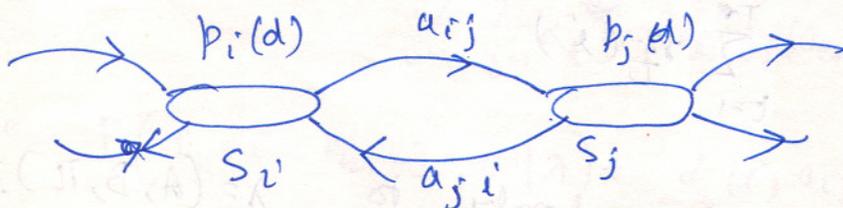
$P_i(d) =$ Prob. of d consecutive observations in state i

$$= (a_{ii})^{d-1} \cdot (1 - a_{ii})$$

\Rightarrow Exponential state duration density model

However, in ~~practice~~ ^{reality}, ~~these~~ most physical systems follow different state duration density model.

Hence the HMM is modified as:



$\&$ $a_{ii} \neq 0$ for all i .

Ex.

in fact:

$$a_{ii} = 0; \quad \sum_{j=1}^N a_{ij} = 1.$$

\rightarrow Transition is made only after the appropriate number of observations have occurred in the state. (as specified by the duration density).

Hence for calculating $p(O|\lambda)$:

$$\alpha_t(i) = p(O_1, O_2, \dots, O_t, S_t \text{ ends at } t | \lambda)$$

Let us assume that a total r states have been visited during the first t observations & let us denote the states as a_1, a_2, \dots, a_r with durations associated with each state as d_1, d_2, \dots, d_r .

$$\therefore a_t = S_t$$

$$\sum_{s=1}^r d_s = t$$

$$\therefore \alpha_t(i) = \sum_{a_1} \sum_d \pi_{a_1} \cdot p_{a_1}(d_1) \cdot p(O_1, O_2, \dots, O_{d_1} | a_1) \cdot a_{a_1 a_2} \cdot p_{a_2}(d_2) \cdot p(O_{d_1+1}, \dots, O_{d_1+d_2} | a_2) \cdot \dots \cdot a_{a_{r-1} a_r} \cdot p_{a_r}(d_r) \cdot p(O_{d_1+d_2+\dots+d_{r-1}+1}, \dots, O_t | a_r)$$

By induction:

$$\alpha_t(j) = \sum_{i=1}^N \sum_{d=1}^D \alpha_{t-d}(i) a_{ij} p_j(d) \prod_{s=t-d+1}^t b_j(O_s)$$

..... (A)

where D is the max. duration within any state.

To initialize the computation of $\alpha_t(j)$ we we,

$$\alpha_1(i) = \pi_i p_i(1) \cdot b_i(O_1)$$

$$\alpha_2(i) = \pi_i p_i(2) \cdot \prod_{s=1}^2 b_i(O_s) + \sum_{\substack{j=1 \\ j \neq i}}^N \alpha_1(j) a_{ji} p_i(1) b_i(O_2)$$

$$\alpha_3(i) = \pi_i p_i(3) \prod_{s=1}^3 b_i(O_s) + \sum_{d=1}^2 \sum_{\substack{j=1 \\ j \neq i}}^N \alpha_{3-d}(j) a_{ji} p_i(d) \prod_{s=4-d}^3 b_i(O_s)$$

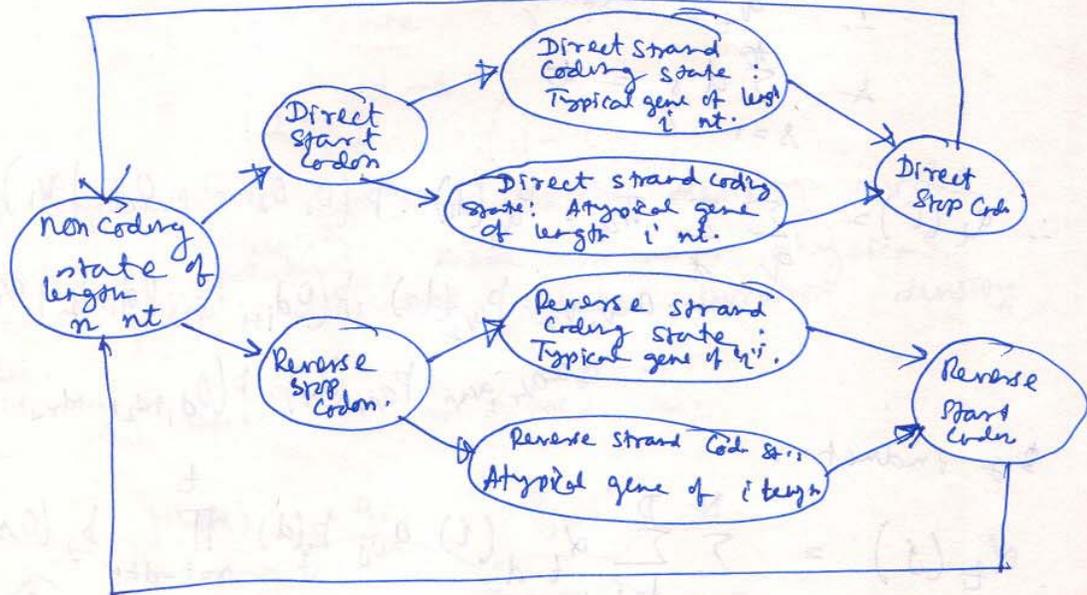
etc. until $\alpha_D(i)$.

Then eqn (A) should be used for all $t > D$.

$$\therefore P(O|A) = \sum_{i=1}^N \alpha_T(i)$$

Ref. "GeneMark.hmm: new solutions for gene finding,"
A. V. Laskov & M. Borodovsky, *Nucleic Acids Research*, 1998, vol. 26, no. 4, 1107-1115.

⊕ For prokaryotic nucleotide sequences



"Typical" genes: Vertically transferred.

"Atypical" " : Horizontally "

Accuracy Measure:

AP = Actual positive (set of genes)

AN = " -ve (" ")
[according to Annotation]

PP = Predicted Positive

PN = " -ve

[following HMM prediction]

$$\begin{aligned} \therefore TP &= \text{True Positive} = \{PP \cap AP\} \\ TN &= \text{" -ve} = \{PN \cap AN\} \\ FP &= \text{False +ve} = \{PP \cap AN\} \\ FN &= \text{" -ve} = \{PN \cap AP\} \end{aligned}$$

$$\begin{aligned} \therefore \text{Sensitivity} = S_n &= \frac{|TP|}{|AP|} \\ \& \text{ Specificity} = S_p &= \frac{|TN|}{|PN|} \end{aligned}$$

$|X| \Rightarrow$
Cardinality
 \uparrow set X.

$$CC = \text{Corr. Coeff.} = \frac{(TP)(TN) - (FP)(FN)}{\sqrt{(|PP|)(|PN|)(|AP|)(|AN|)}}$$

AC = Approximate Cor.

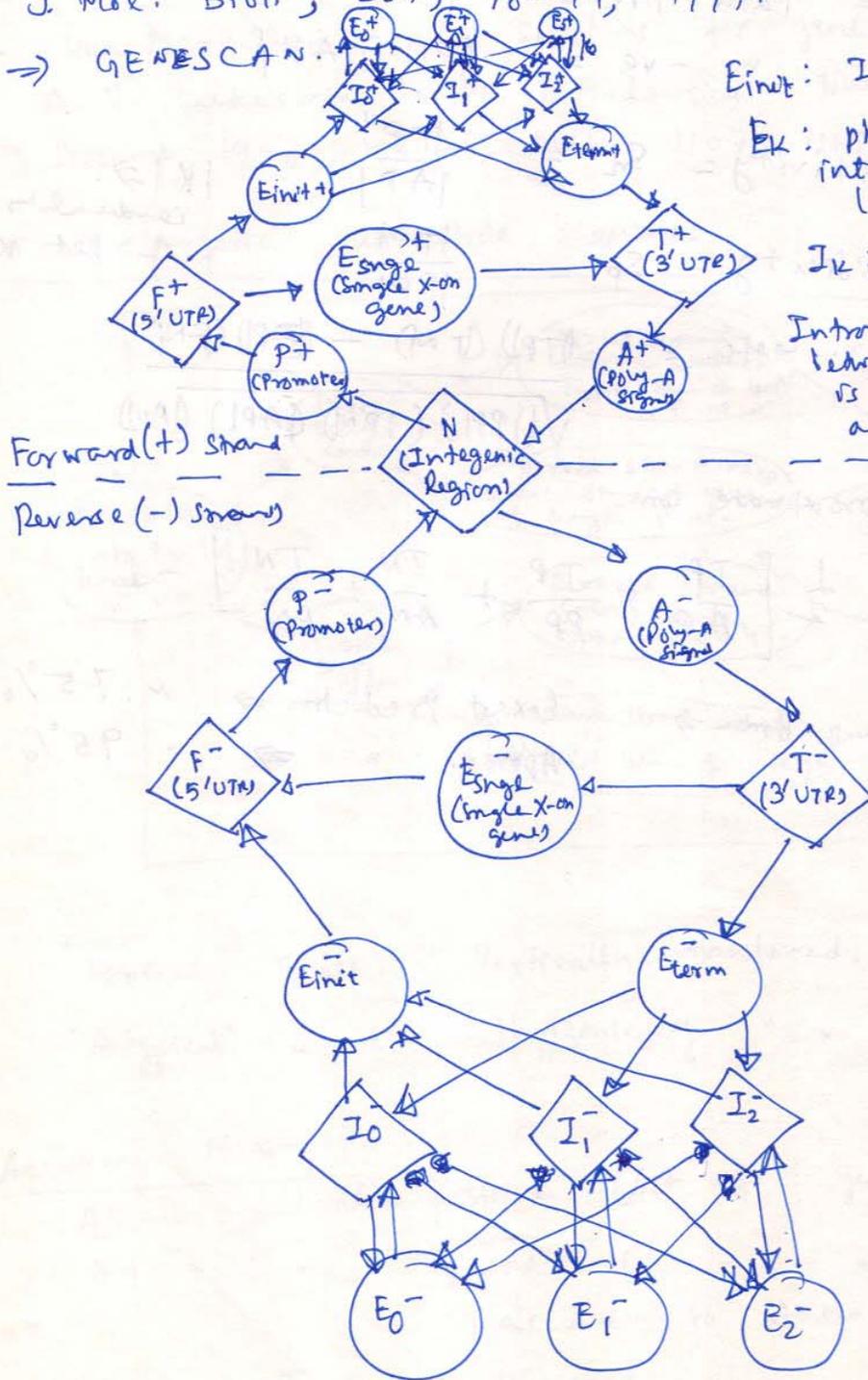
$$\approx \frac{1}{2} \left[\frac{TP}{AP} + \frac{TP}{PP} + \frac{TN}{AN} + \frac{TN}{PN} \right] - 1$$

Gene Mark. Ann \Rightarrow

Exact Prediction $\Rightarrow \sim 75\%$
Approx. " $\Rightarrow \sim 95\%$

Ref.: "Prediction of Complete Gene Structures in Human Genomic DNA", C. Burge & S. Karlin, J. Mol. Biol., 267, 78-94, 1997.

⊕ ⇒ GENESCAN.



Einit: Initial Exon

Ek: phase k internal Exon (k=0,1,2).

Ik: phase k intron

Introns falls between codons
is at phase 0, after first base
of codon at phase 1
& after 2nd. base of a codon
at phase 2.

Internal exons (Ek) are similarly denoted by the phase of its precedes Intron.

Forward (+) strand
Reverse (-) strand

Performance 75-80% exon correctly detected.

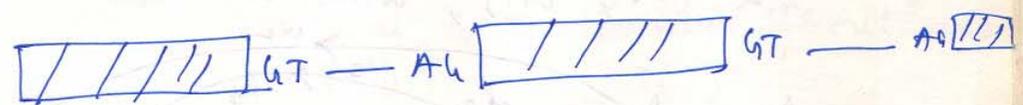
GENSCAN
Gene Xon Performance:

cc: 0.93 $S_n = 0.95 - 0.98$ $S_p = 0.90 - 0.94$
exons correct: 0.76 - 0.79 ; exons overlapped: 0.91 - 0.96

② Similarity-Based Approaches: (Should be discussed with the topic on sequence alignment)

Prob. defn Given a known target protein and a genomic sequence, find a set of substrings (candidate exons) of the genomic sequence whose concatenation (splicing) best fits the target.

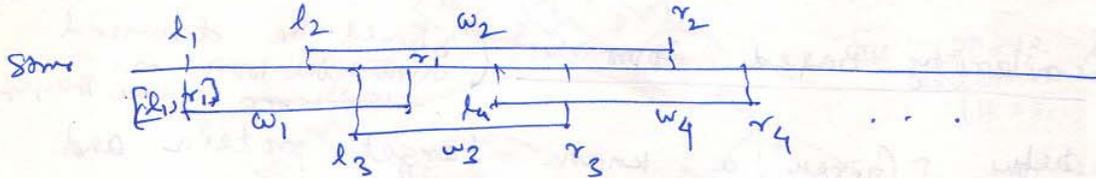
- Find all local similarities between the genomic sequence & the target protein sequence.
- Each substring from the genomic sequence that exhibits sufficient similarity to the target protein could be considered as a putative exon.
- The putative exons so chosen may lack the canonical exon-flanking dinucleotides AG and GT, but we can extend or shorten them slightly to make sure they are flanked by AG & GT:



- The resulting candidate may contain overlapping substrings, out of which best possible non-overlapping substrings to be chosen with a maximal match to the target protein.

• Exon - Chain Problem

Weighted Interval: $[l, r] \Rightarrow$ exon & weight w .
 $\Rightarrow (l, r, w), w > 0$.



Chain: any set of non-overlapping intervals

weight of a chain: Sum of ^{interval} weights (intervals)

Maximum Chain: ^{chain} Maximum weight among all possible chains.

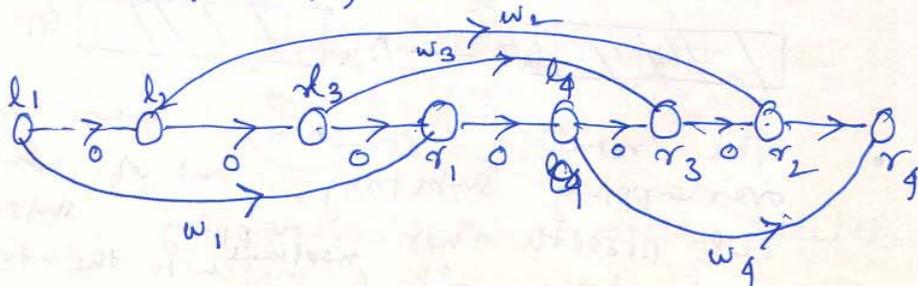
Problem Defn:

Given a set of putative exons, find a max. set of nonoverlapping putative exons.

Soln: $G(V, E): V = 2n$ vertices, n of which represent left positions & n of which represent right position (let them be distinct).

Let $V = (v_1, v_2, \dots, v_{2n})$ be the ordered array of vertices.

$E: \langle l_i, r_i \rangle$ with w_i
 & (v_i, v_{i+1}) , with 0 weight.

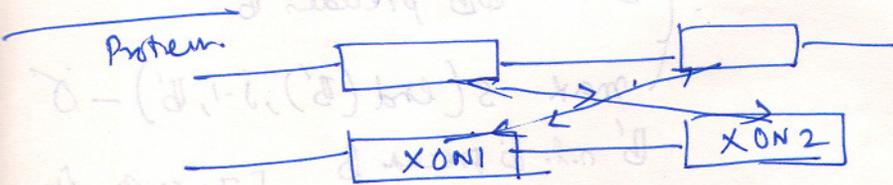


The objective ~~problem~~ is to find a path from v_1 to v_n with maximal weight.

Use of dynamic prog.

At vertex v_i , ~~weight~~ $s_i = \max \{ s_j + w, s_{i-1} \}$
 s.t. v_j is connected to v_i .
 $\Leftrightarrow \exists$ an interval $[v_j, v_i]$ with w

Shortcoming:



An infeasible chain that might have a maximal score.

o Spliced Alignment:

Problem Statement:

Given Genomic Seq. G , target seq. T and a set of candidate exon (blocks) B ,

compute a chain of candidate exons Γ s.t. the global alignment score $s(\Gamma^*, T)$ is maximum among all possible chains of B .

Solu.: Compute,
 1. Score (v, j, B) , $b \in B$, & $G(v) \in b$.
 so that $\rightarrow b$ & maximal score terminating at $G(v)$ for with target seq. $T(0) T(1) \dots T(j)$
 using DP.

2. Finally, Let $end(b)$ denote the pos end position in G . & m be the length of target seq. T .

∴ maximal spaced score = $\max_B s(\text{end}(B), m, B)$

Recurrence relation.

Let i be the starting vertex of B .

$$s(i, j, B) = \max \begin{cases} s(i, j-1, B) - \sigma & \text{[deletion in the target seq]} \\ \max_{B' \text{ n.t. } B' \text{ precedes } B} s(\text{end}(B'), j-1, B') + \delta(g_i, t_j) \\ \max_{B' \text{ n.t. } B' \text{ precedes } B} s(\text{end}(B'), j-1, B') - \sigma & \text{[Insertion in genome seq]} \end{cases}$$

Let i' be an intermediate or end position for B .

$$s(i, j, B) = \max \begin{cases} s(i, j-1, B) - \sigma & \text{! deletion} \\ \max_{B' \text{ precedes } B} (s(\text{end}(B'), j-1, B') + \delta(g_i, t_j)) \\ s(i, j-1, B) - \sigma & \text{/ insertion} \\ s(i-1, j-1, B) + \delta(g_i, t_j) \end{cases}$$

Small m Exon blocks should be filtered out, so that missalng effect causing incorrect predictions is avoided.

e.g. "filtering of" \Rightarrow could be matched with ^{easy} any alphabet longer for any random sequence of the alphabet.