

## Week 14 Annotation

Monday, 18 April 2022 8:40 AM

### ① Computation phase

- ↳ proteins are identified
- ↳ ab initio or evidence based

### ② annotation phase

- ↳ synthesize annotation

#### STEP 1 Repeat identification.

- ↳ transposons, LINEs, SINEs, viral seqs.
- ↳ RepeatMasker

#### STEP 2 Evidence alignment

eg. 1 RNA seq. aligned to assembly

eg. 2. ab initio gene prediction

(eg. motifs → AUG (start),  
TATAA boxes, etc.)

stop codons,

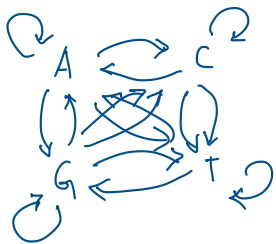
#### STEP 3: UniProtKB, SwissProt, PDB.

(linking genes to proteins)

#### STEP 4: visualization (genome browsers)

#### HMM's

$\{A, C, G, T\}$



TRANSITION

$$PROB = a_{st}$$

= prob. of  
transitioning  
from state 's' to  
state 't'

$$P(X_i = t \mid X_{i-1} = s) \quad \text{①}$$

$$a_{st} = \underbrace{\quad \quad \quad}$$

$$\begin{aligned} p(x) &= P(x_L, x_{L-1}, x_{L-2}, \dots, x_1) \\ &= P(x_L | x_{L-1}, x_{L-2}, \dots, x_1) \times P(x_{L-1} | x_{L-2}, \dots, x_1) \\ &\quad \times \dots \times P(x_2 | x_1) \times P(x_1) \\ &= P(x_L | x_{L-1}) \times P(x_{L-1} | x_{L-2}) \times P(x_{L-2} | x_{L-3}) \\ &\quad \times \dots \times P(x_2 | x_1) \times \underline{P(x_1)} \end{aligned}$$

Markov property →

$$\therefore p(x) = P(x_1) \prod_{i=2}^L \overset{L}{\underset{\text{observed}}{a_{x_{i-1}} x_i}} \quad (2)$$

in your HMM:

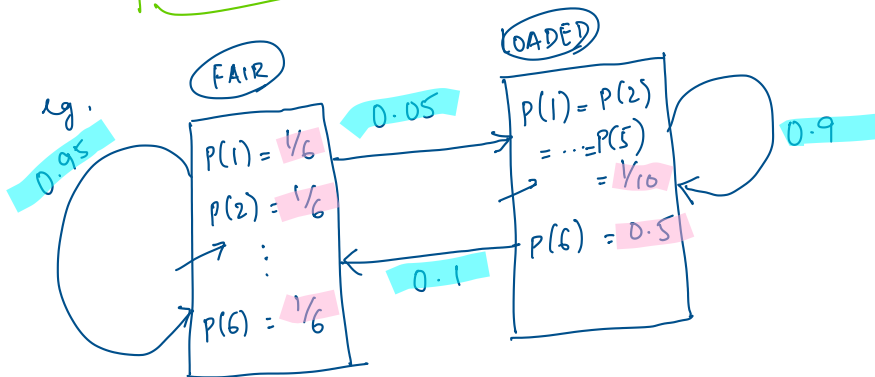
$$\begin{aligned} &\boxed{\text{state sequence path} = \pi} \\ a_{kl} &= P(\pi_i = l | \pi_{i-1} = k) \end{aligned}$$

TRANSITION PROBABILITIES

'SYMBOLS' → (LABELS) → hidden states

$$e_k(b) = P(\underbrace{x_i = b}_{\text{EMISION PROBABILITIES}} | \pi_i = k)$$

$x$  = unobserved label path



VITERBI ALGORITHM

to pick best 'next' step, compute the

... probability

seq.  
n, π

hyper...  
 $\pi^* = \arg \max_{\pi} P(x, \pi)$   $\uparrow$  state.

$v_k(i)$  = prob. of the most probable path ending in state  $k$  with observed symbol  $i$

then

$$v_l(i+1) = e_l(x_{i+1}) \max_k (v_k(i) a_{kl})$$

$\downarrow$  Prob. of observing  $l$  in  $i^{th}$  state.       $\downarrow$  transition of  $k \rightarrow l$

eg. 2 1 2 3 6 2 4 ...  
 1 2 3 4 5 6 7

$$v_L(6,5) = e_L(6) \left( \max_{\underline{L}} (v_{\underline{L}}(3,4), v_{\underline{F}}(3,4)) \right)$$

Fair coin or a biased coin

FAIR COIN

	H	T
H	0.5	0.5
T	0.5	0.5

LOADED COIN

	H	T
H	0.7	0.3
T	0.7	0.3

EMISSION PROBS  $\nearrow$

TRANSITION PROBS.

	F	L
F	0.95	0.05
L	0.1	0.90

FORWARD ALGORITHM

Full prob. of a sequence  $x$

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

$$f_k^{(i)} = P(x_1, x_2, \dots, x_i, \pi_i = k)$$

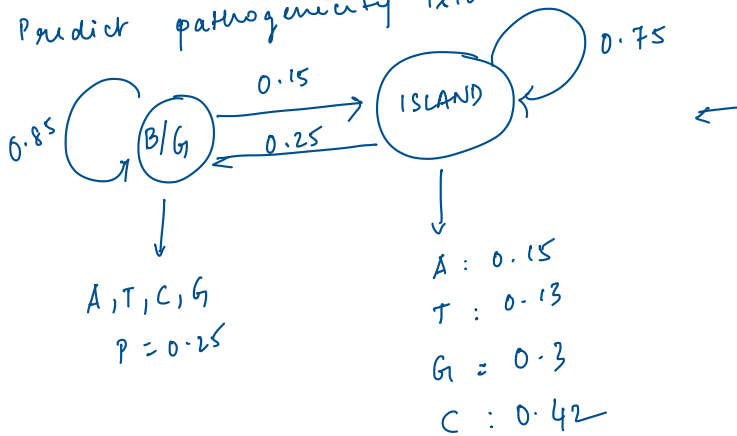
Prob. of the observed sequence up to and including  $x_i$ ; such that  $\pi_i = k$

$$f_k^{(i+1)} = e_l^{(x_{i+1})} \sum_k f_k^{(i)} a_{kl}$$

$$G(i|x) = \sum_k P(\pi_i = k | x) g(k) \quad \leftarrow$$

$g(k) = \text{PRIOR on states}$

Predict pathogenicity islands.



TRANSITION

	B	I
B	0.85	0.15
I	0.25	0.75

EMISSION

	A	T	C	G
B	0.25	0.25	0.25	0.25
I	0.15	0.13	0.3	0.42

