

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,

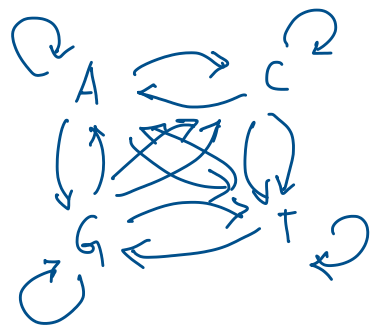
STEP 3 : UniProtKB, SwissProt, PDB.

(linking genes to proteins)

STEP 4 : visualization (genome browsers)

HMM's

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



TRANSITION

$$\text{PROB} = a_{st}$$

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

$$a_{st} = \underbrace{P(X_i = t \mid X_{i-1} = s)} \leftarrow (1)$$

$$\begin{aligned} p(x) &= P(X_L, X_{L-1}, X_{L-2}, \dots, X_1) \\ &= P(X_L \mid X_{L-1}, X_{L-2}, \dots, X_1) \times P(X_{L-1} \mid X_{L-2}, \dots, X_1) \\ &\quad \times \dots \times P(X_2 \mid X_1) \times P(X_1) \\ &= P(X_L \mid X_{L-1}) \times P(X_{L-1} \mid X_{L-2}) \times P(X_{L-2} \mid X_{L-3}) \\ &\quad \times \dots \times P(X_2 \mid X_1) \times \underline{P(X_1)} \end{aligned}$$

Markov
property \rightarrow

$$\therefore p(x) = P(X_1) \prod_{i=2}^L a_{x_{i-1} x_i} \leftarrow (2)$$

observed

in your HMM:

state sequence path = π \leftarrow

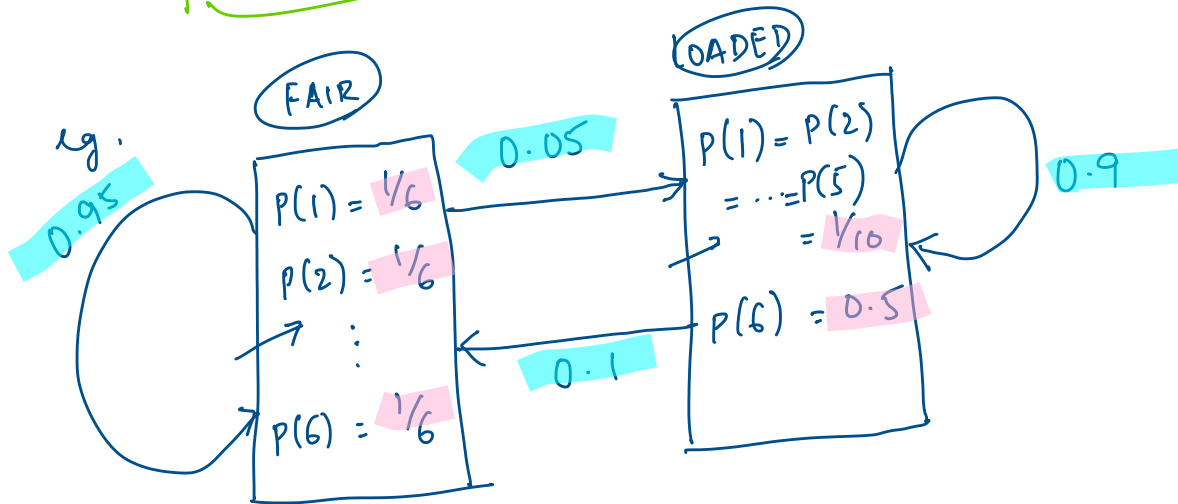
$$a_{kl} = P(\pi_i = l \mid \pi_{i-1} = k) \leftarrow$$

TRANSITION
PROBABILITIES

'SYMBOLS' \rightarrow (LABELS) \rightarrow hidden states

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

x = unobserved label path



VITERBI ALGORITHM

to pick best 'next' step, compute the highest probability

$$\pi^* = \arg \max_{\pi} P(x, \pi)$$

\downarrow seq.
 \uparrow state.

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

then

$$v_i(i+1) = e_{\ell_i}(x_{i+1}) \max_k (v_k(i) a_{k\ell_i})$$

\downarrow
 transition

l

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

eg.
 $\begin{matrix} & \downarrow & \downarrow & \downarrow & \downarrow & & \\ 2 & 1 & 2 & 3 & 6 & 2 & 4 & \dots \end{matrix}$
 $\begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \end{matrix}$

$$V_L(6,5) = e_L(6) \max(V_L(3,4), V_F(3,4))$$

Fair coin on 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 PROB.

	F	L
F	0.95	0.05
L	0.1	0.90