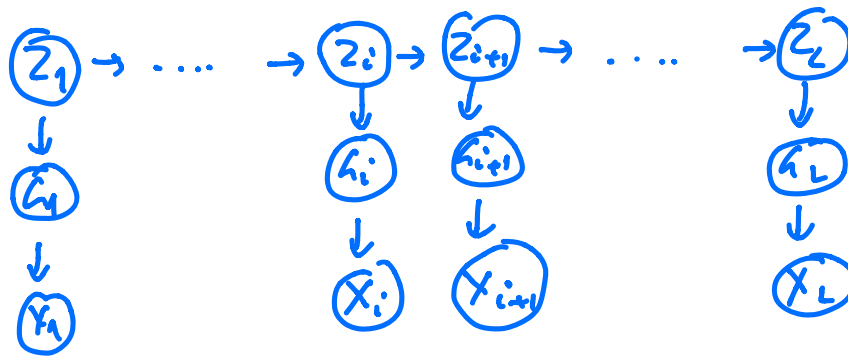


W07 - Hidden Markov Models - Inference  
Expectation - Maximization

Inference by dynamic programming



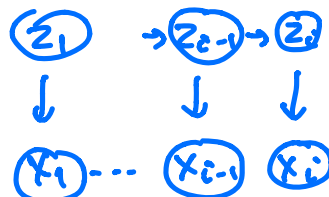
Data =  $[Y_1 \dots Y_L]$  reads  $X_i = \{n_a^i, n_c^i, n_g^i, n_t^i\}$

Inference  $[Z_1 \dots Z_L]$  ancestry  $Z_i \in \{A, B, A\}$

Ignored  $[G_1 \dots G_L]$  genome (interpreted)  $G_i \in \{a, c, g, \dots\}$

Forward algorithm

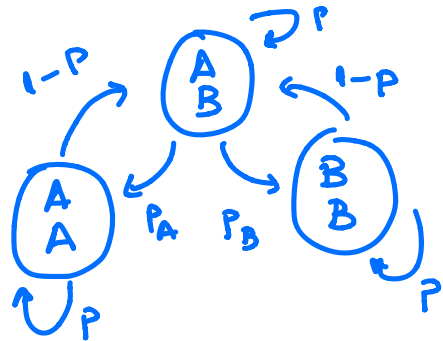
$$f_z(i) = P(Y_1 \dots X_{i-1} X_i Z_i = z)$$



$$f_z(i) = \sum_{Z_{i-1} \in \{A, B, A\}} f_{Z_{i-1}}(i-1) \cdot P(z | Z_{i-1}) P(X_i | z)$$

Forward

$$f_z(i) = \sum_{z_{i-1} = \begin{matrix} ABA \\ ABB \end{matrix}} f_{z_{i-1}}(i-1) P(z|z_{i-1}) P(x_i|z)$$



$$f_z(i) = f_{AA}(i-1) P(z|AA) P(x_i|z) \\ + f_{BB}(i-1) P(z|BB) P(x_i|z) \\ + f_{AB}(i-1) P(z|AB) P(x_i|z)$$

z=AA

$$f_{AA}(i) = f_{AA}(i-1) P \cdot P(x_i|AA) \\ + f_{AB}(i-1) P_A P(x_i|AA)$$

z=BB

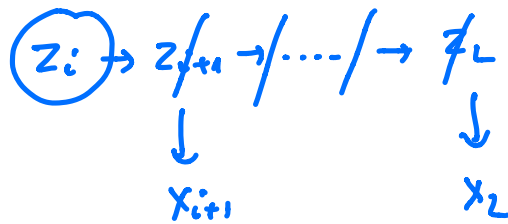
$$f_{BB}(i) = f_{BB}(i-1) P P(x_i|BB) \\ + f_{AB}(i-1) P_B P(x_i|BB)$$

z=AB

$$f_{AB}(i) = f_{AB}(i-1) P P(x_i|AB) \\ + f_{AA}(i-1) (1-P) P(x_i|AB) \\ + f_{BB}(i-1) (1-P) P(x_i|AB)$$

Int:  $f_z(i) = P_r(z) P(x_i|z)$

## Backward Algorithm



$$b_2(i) = P(x_{i+1} \dots x_L | z_i = z)$$

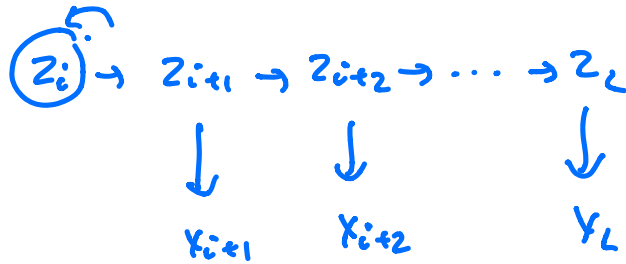
$$b_2(i) = \sum_{z_{i+1}} \dots \sum_{z_L} P(x_{i+1} z_{i+1} \dots x_L z_L | z_i = z)$$

$$= \sum_{z_{i+1}} \sum_{z_{i+2}} \dots \sum_{z_L} P(x_{i+1} z_{i+1} | z_i = z) \cdot P(x_{i+2} z_{i+2} \dots x_L z_L | z_{i+1})$$

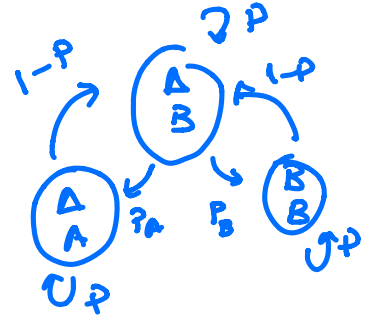
$$= \sum_{z_{i+1}} P(x_{i+1} z_{i+1} | z_i = z) b_{2, z_{i+1}}^{(i+1)}$$

$$= \sum_{z_{i+1}} P(x_{i+1} | z_{i+1}) P(z_{i+1} | z) \cdot b_{2, z_{i+1}}^{(i+1)}$$

$$b_2(i) = \sum_{z_{i+1}} b_{2, z_{i+1}}^{(i+1)} \cdot P(z_{i+1} | z) P(x_{i+1} | z_{i+1})$$



$$\begin{aligned}
 b_2(i) &= b_{AA}(i+1) P(AA|Z) P(Y_{i+1}|AA) \\
 &+ b_{BB}(i+1) P(BB|Z) P(Y_{i+1}|BB) \\
 &+ b_{AB}(i+1) P(AB|Z) P(Y_{i+1}|AB)
 \end{aligned}$$



$$\begin{aligned}
 \underline{Z=AA} \\
 b_{AA}(i) &= b_{AA}(i+1) p P(Y_{i+1}|AA) \\
 &+ b_{AB}(i+1) (1-p) P(Y_{i+1}|AB)
 \end{aligned}$$

$$\begin{aligned}
 \underline{Z=BB} \\
 b_{BB}(i) &= b_{BB}(i+1) p P(Y_{i+1}|BB) \\
 &+ b_{AB}(i+1) (1-p) P(Y_{i+1}|AB)
 \end{aligned}$$

$$\begin{aligned}
 \underline{Z=AB} \\
 b_{AB}(i) &= b_{AB}(i+1) p P(Y_{i+1}|AB) \\
 &+ b_{AA}(i+1) p_A P(Y_{i+1}|AA) \\
 &+ b_{BB}(i+1) p_B P(Y_{i+1}|BB)
 \end{aligned}$$

Initialization

$$b_{AA}(L) = b_{AB}(L) = b_{BB}(L) = 1$$

$$P(Y_1 \dots Y_L | Z_L = AA) = 1$$

$P(x_1 \dots x_L) = \text{Likelihood}$

$$\begin{aligned} P(x_1 \dots x_L) &= \sum_{z_i} P(\overbrace{x_1 \dots x_{i-1} y_i z_i}^B \overbrace{x_{i+1} \dots x_L}^A) \\ &= \sum_{z_i} P(x_{i+1} \dots x_L | z_i) \\ &\quad P(x_1 \dots x_{i-1} y_i z_i) \\ &= \sum_{z_i} b_{z_i}(i) \cdot f_{z_i}(i) \end{aligned}$$

$$P(x_1 \dots x_L) = \sum_{z_i} f_{AA}(i) b_{AA}(i) + f_{BB}(i) b_{BB}(i) + f_{AB}(i) b_{AB}(i)$$

Good test for forward/backward correctness

Decoding - Finally the Inference we are up to!

$$P(z_i | y_1 \dots y_L) = P(z_i | \text{Data}, \text{HMM } P)$$

3 cases

$$P(z_i = AA | y_1 \dots y_L)$$

$$P(z_i = BB | y_1 \dots y_L)$$

$$P(z_i = AB | y_1 \dots y_L)$$

$$\frac{\quad}{\quad} = 1$$

$$P(z_i | y_1 \dots y_L) = \frac{P(y_1 \dots x_i z_i \dots y_L)}{P(y_1 \dots y_L)}$$

$$= \frac{P(\overset{B}{y_1 \dots x_i} z_i \overset{A}{x_{i+1} \dots y_L})}{P(y_1 \dots y_L)}$$

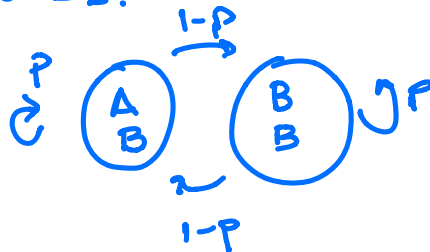
$$= \frac{P(x_{i+1} \dots y_L | z_i) \cdot P(y_1 \dots x_i | z_i)}{P(y_1 \dots y_L)}$$

$$P(z_i | y_1 \dots y_L) = \frac{b_{z_i}(i) \cdot f_{z_i}(i)}{\sum_z b_z(i) f_z(i)}$$

→ class code

W07 - How much does the value of  $p$  matter?

Consider only backcrosses to BB. AA not possible



$$P(z_i | z_{i-1}) = \begin{cases} p & z_i = z_{i-1} \\ 1-p & z_i \neq z_{i-1} \end{cases}$$

$$P\left(\begin{matrix} A \\ B \end{matrix} \dots \begin{matrix} A \\ B \end{matrix} \mid \begin{matrix} A \\ B \end{matrix} \begin{matrix} A \\ B \end{matrix}\right) = p^n (1-p)$$

$$P\left(\begin{matrix} B \\ B \end{matrix} \dots \begin{matrix} B \\ B \end{matrix} \mid \begin{matrix} B \\ B \end{matrix} \begin{matrix} B \\ B \end{matrix}\right) = p^n (1-p)$$

$$\langle n \rangle = \frac{p}{1-p}$$

*Drosophila*  $\langle n \rangle = L/2$

chr 2L len =  $23.01 \cdot 10^6$

chr 4 len =  $1.35 \cdot 10^6$

$p = 0.99999991$  ?

$p = 0.99999852$  .

does it matter?

→ class code

## Maximum Likelihood estimates

To do ML estimates, we need "labelled data"

i.e. Data for which  $z_i$  are known.

$$D = \left\{ \left\{ X_1^m z_1^m \dots X_L^m z_L^m \right\}_{m=1}^M \right\} \quad M = \# \text{ of flies}$$

$$P(D|HMM) = \prod_{m=1}^M P(X_1^m z_1^m \dots X_L^m z_L^m | P)$$

$$= \prod_{m=1}^M P(z_1^m) P(X_1^m | z_1^m) \prod_{i=2}^L P(z_i^m | z_{i-1}^m) P(X_i^m | z_i^m)$$

$$C_i^m(AA \rightarrow AA) + C_i^m(AB \rightarrow AB) = C_i^m(\text{same})$$

$$C_i^m(AB \rightarrow BB) + C_i^m(BB \rightarrow AB) = C_i^m(\text{break})$$

$$P(D|P) = \prod_{m=1}^M P(z_1^m) \prod_{i=1}^L P^{C_i^m(\text{same})} (1-P)^{C_i^m(\text{break})} \cdot \prod_{i=1}^L P(X_i^m | z_i^m)$$

$$\log P(D|P) \propto C \sum_{m=1}^M \sum_{i=1}^L C_m^i(s) \log P + \sum_m \sum_i C_m^i(b) \log(1-P)$$



$$\log P(D|P) \propto C(s) \log P + C(b) \log(1-P)$$

$$C(s) = \sum_m \sum_i [C_m^i(AB \rightarrow AB) + C_m^i(BB \rightarrow BB)]$$

$$C(b) = \sum_m \sum_i [C_m^i(AB \rightarrow BB) + C_m^i(BB \rightarrow AB)]$$

$$\frac{\delta \log P(D|P)}{\delta P} = \frac{C(s)}{P} - \frac{C(b)}{1-P} = 0$$

$$\frac{C(s)}{P_{ML}} = \frac{C(b)}{1-P^*}$$

$$P^* = \frac{C(s)}{C(s) + C(b)}$$