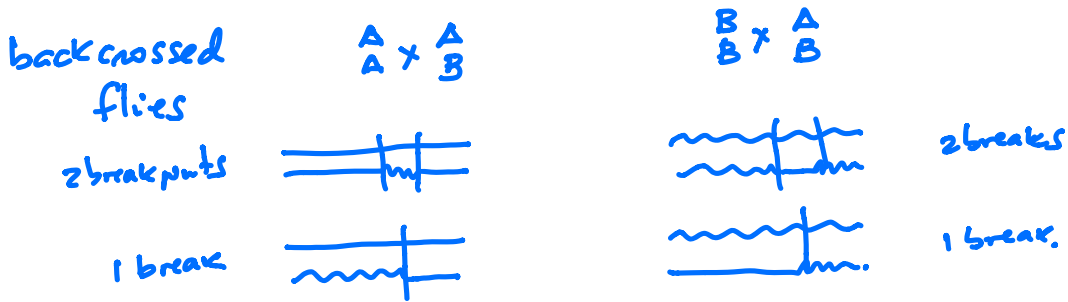
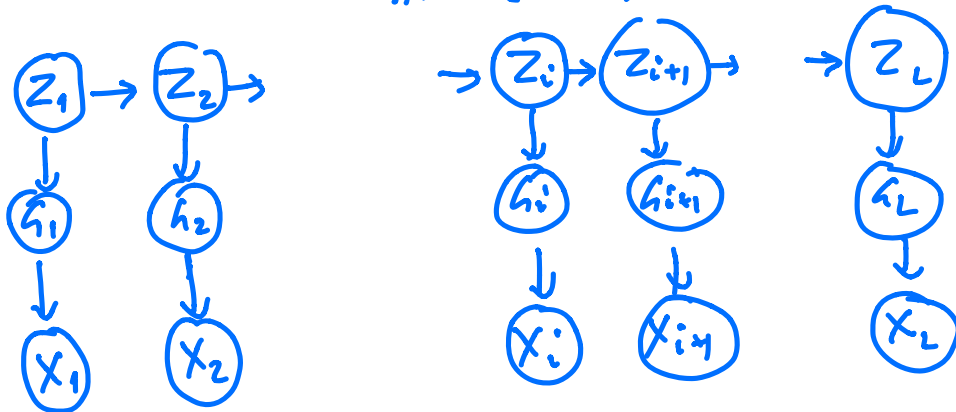


W06 - Ancestry Inference from RAA-seq reads
using a hidden Markov model (HMM)



$Z = \left\{ \begin{matrix} A & B & A \\ A & B & B \end{matrix} \right\}$ ancestry, Z_i for each position i

HMM (DAG)



$h_i =$ the znts at the genome in that position. (A A A A C C
a, c, g, t, c, g, ...))

$X_i = (n_a^i \ n_c^i \ n_g^i \ n_t^i)$ the reads at position i

$$\begin{aligned}
 P(x_1, g_1, z_1, \dots, x_L, g_L, z_L) &= P(x_L, g_L, z_L \mid x_{L-1}, g_{L-1}, z_{L-1}) \\
 &\quad \text{from DAG} \quad \curvearrowright \quad : \\
 &\quad \quad \quad P(x_2, g_2, z_2 \mid x_1, g_1, z_1) \\
 &\quad \quad \quad P(x_1, g_1, z_1) \\
 &= P(x_1, g_1, z_1) \cdot \prod_{i=2}^L P(x_i, g_i, z_i \mid x_{i-1}, g_{i-1}, z_{i-1})
 \end{aligned}$$

From DAG:

$$\begin{aligned}
 P(x_i, g_i, z_i \mid x_{i-1}, g_{i-1}, z_{i-1}) &= P(z_i \mid z_{i-1}) \\
 &\quad P(g_i \mid z_i) \\
 &\quad P(x_i \mid g_i)
 \end{aligned}$$

} x_i } data

} z_i } variables we want to know about

} g_i } hidden variables, we need them to go from $z_i \rightarrow x_i$

but their actual values are not interesting to us

$$P(x_1, z_1, \dots, x_L, z_L) = \sum_{g_1} \dots \sum_{g_L} P(x_1, g_1, z_1, \dots, x_L, g_L, z_L)$$

$$= P(z_1) \sum_{g_1} P(x_1 \mid g_1) \cdot P(g_1 \mid z_1)$$

$$\prod_{i=2}^L P(z_i \mid z_{i-1}) \cdot \underbrace{\sum_{g_i} P(x_i \mid g_i) P(g_i \mid z_i)}_{P(x_i \mid g_i)}$$

$$P(x_1 z_1 \dots x_L z_L) = P(z_1) P(x_1 | z_1)$$

$$\prod_{i=2}^L P(z_i | z_{i-1}) P(x_i | z_i)$$

The parameters, where are they?

$P(z_i | z_{i-1})$ " 3 probability distributions

$$P(A|A) = P$$

$$P(A|B) = 0$$

$$P(A|B) = P_A$$

$$P(B|A) = 0$$

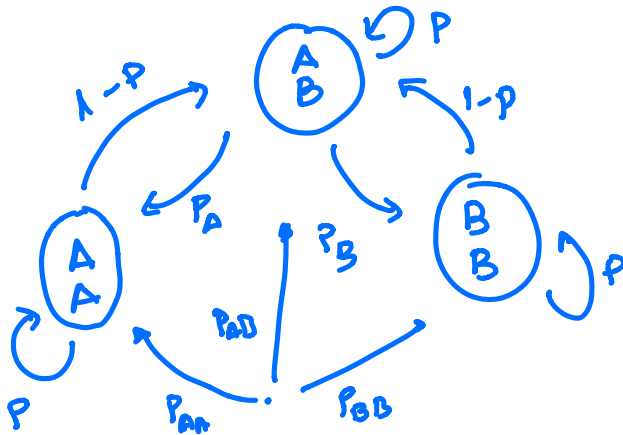
$$P(B|B) = P$$

$$P(B|B) = P_B$$

$$P(A|B) = 1-P$$

$$P(B|A) = 1-P$$

$$P(B|A) = P$$



$$P + P_A + P_B = 1$$

State diagram \neq DAG

P? $\frac{-n-}{m-m}$ length between breaks = geometric distribution
 $P^n (1-P) = P(n) \quad \langle n \rangle = \frac{P}{1-P}$

the "emission" parameters



$$P(G_i | Z_i) \begin{cases} P(G_i | AA) = \begin{cases} 1 & \text{if } G_i = g_i^A g_i^A \\ 0 & \text{otherwise} \end{cases} \\ P(G_i | BB) = \begin{cases} 1 & \text{if } G_i = g_i^B g_i^B \\ 0 & \text{otherwise} \end{cases} \\ P(G_i | AB) = \begin{cases} 1 & \text{if } G_i = g_i^A g_i^B \\ 0 & \text{otherwise} \end{cases} \end{cases}$$

$$P(X_i | G_i)$$

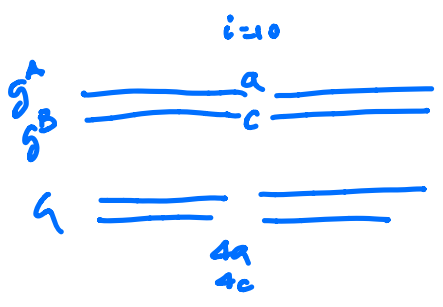
- no parameters
- assumes no errors in the genome
- Andolfato has an ϵ error parameter

$$P(X_i | G_i) = \begin{matrix} n_a & n_c+n_g+n_t \\ P_1 & q_1 \end{matrix} \left. \begin{matrix} G_i = aa \\ \vdots \end{matrix} \right\} 4 \text{ cases}$$

$$P(X_i | G_i) = \begin{matrix} n_a+n_c & n_g+n_t \\ P_2 & q_2 \end{matrix} \left. \begin{matrix} G_i = ac \\ \vdots \end{matrix} \right\} 6 \text{ cases}$$

$$\left. \begin{matrix} P_1 = 1 - \epsilon \\ q_1 = \epsilon/3 \end{matrix} \right\} \begin{matrix} P_2 = \frac{1-\epsilon}{2} \\ q_2 = \epsilon/2 \end{matrix}$$

$$P_1 + 3q_1 = 1 \quad P_2 + P_2 + q_2 + q_2 = 1$$



$$\begin{aligned} P(G_{10} | AA) &= 1 \text{ if } G_{10} = aa \\ P(G_{10} | BB) &= 1 \text{ if } G_{10} = cc \\ P(G_{10} | AB) &= 1 \text{ if } G_{10} = ac \end{aligned}$$

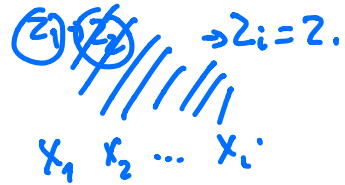
$$P(X_{10} | AA) = \sum_{\zeta_{10}} P(X_{10} | \zeta_{10}) P(\zeta_{10} | AA)$$

$$= P(X_{10} | aa) = \binom{4+4}{1} \binom{0}{1} = \binom{8}{1}$$

$$P(X_{10} | AB) = P(X_{10} | ac) = \binom{8}{2} \binom{0}{2} = \binom{8}{2}$$

$$P(X_{10} | BB) = P(X_{10} | cc) = \binom{4}{1} \binom{4}{1}$$

The forward algorithm



$$f_2(i) = P(x_1 x_2 \dots x_i z_i = z)$$

$$= \sum_{z_1} \dots \sum_{z_{i-1}} P(x_1 z_1, x_2 z_2, \dots, x_i z_i = z)$$

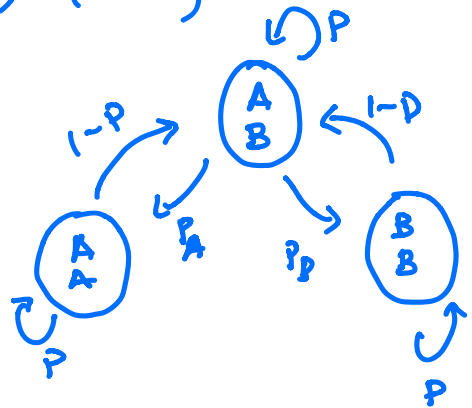
$$= \sum_{z_1} \sum_{z_{i-1}} P(x_1 z_1) \cdot P(x_2 z_2 | x_1 z_1) \dots P(x_{i-1} z_{i-1} | x_{i-2} z_{i-2}) \\ P(x_i z_i = z | x_{i-1} z_{i-1})$$

$$= \sum_{z_{i-1}} \underbrace{\sum_{z_1} \dots \sum_{z_{i-2}} P(x_1 z_1) P(x_2 z_2 | x_1 z_1) \dots P(x_{i-1} z_{i-1} | x_{i-2} z_{i-2})}_{f_{z_{i-1}}(i-1)} \cdot P(x_i z_i = z | x_{i-1} z_{i-1})$$

$$= \sum_{z_{i-1}} f_{z_{i-1}}(i-1) \cdot P(x_i z_i = z | x_{i-1} z_{i-1})$$

$$= \sum_{z_{i-1}} f_{z_{i-1}}(i-1) \cdot P(z_i = z | z_{i-1}) P(x_i | z)$$

$$f_z(i) = \sum_{z_{i-1}} f_{z_{i-1}}(i-1) P(z|z_{i-1}) P(x_i|z)$$



$z=AA$ $z_{i-1}=AA, BB, AB$

$$\begin{aligned} f_{AA}(i) &= + f_{AA}(i-1) P(AA|AA) P(x_i|AA) \\ &+ f_{BB}(i-1) P(AA|BB) P(x_i|AA) \\ &+ f_{AB}(i-1) P(AA|AB) P(x_i|AA) \\ &= f_{AA}(i-1) \cdot p \cdot P(x_i|AA) \end{aligned}$$

+ 0

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

$$\begin{aligned} f_{BB}(i) &= + f_{AA}(i-1) P(BB|AA) P(x_i|AA) \\ &+ f_{BB}(i-1) P(BB|BB) P(x_i|BB) \\ &+ f_{AB}(i-1) P(BB|AB) P(x_i|BB) \\ &= f_{BB}(i-1) \cdot p P(x_i|BB) + f_{AB}(i-1) p_B P(x_i|BB) \end{aligned}$$

$$\begin{aligned}
f_{AB}(i) &= f_{AA}(i-1) P(AB|AA) P(x_i|AB) \\
&\quad + f_{BB}(i-1) P(AB|BB) P(x_i|AB) \\
&\quad + f_{AB}(i-1) P(AB|AB) P(x_i|AB) \\
&= f_{AA}(i-1) \cdot (1-p) P(x_i|AB) \\
&\quad + f_{BB}(i-1) (1-p) P(x_i|AB) \\
&\quad + f_{AB}(i-1) p P(x_i|AB)
\end{aligned}$$

$$f_{AA}(1) = P_r(AA) \cdot P(x_1|AA)$$

$$f_{BB}(1) = P_r(BB) \cdot P(x_1|BB)$$

$$f_{AB}(1) = P_r(AB) \cdot P(x_1|AB)$$

$$P(x_1 \dots x_L) = \sum_{z_i} P(\overbrace{x_1 \dots x_i}^B \mid z_i = z) P(\overbrace{x_{i+1} \dots x_L}^A \mid z_i = z)$$

$$= \sum_{z_i} \frac{P(x_{i+1} \dots x_L \mid x_1 \dots x_i, z_i = z)}{P(x_1 \dots x_i \mid z_i)}$$

$$= \sum_{z_i} \underbrace{f_{z_i}(i) \cdot b_{z_i}(i)}_{\text{wavy line}} \quad \forall i$$

good def of jor algorithm. $= f_{AA}(i) g_{AA}(i) + f_{BB}(i) g_{BB}(i) + f_{AB}(i) g_{AB}(i)$

Backward algorithm

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



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

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

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

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

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

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

$$\begin{aligned} b_{AA}(i) &= b_{AA}(i+1) P(AA|AA) P(x_{i+1} | AA) \\ &\quad + b_{BB}(i+1) P(BB|AA) P(x_{i+1} | BB) \\ &\quad + b_{AB}(i+1) P(AB|AA) P(x_{i+1} | AB) \\ &= p b_{AA}(i+1) P(x_{i+1} | AA) \\ &\quad + (1-p) b_{AB}(i+1) P(x_{i+1} | AB) \end{aligned}$$

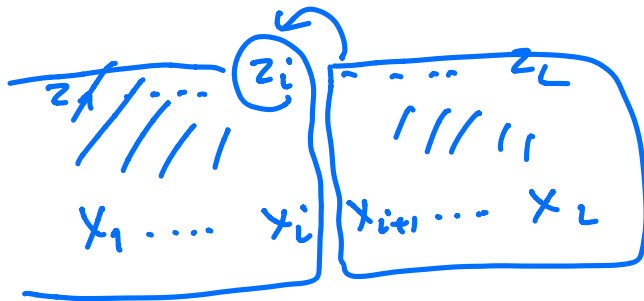
$$\begin{aligned} b_{BB}(i) &= + b_{AA}(i+1) P(AA|BB) P(x_{i+1} | AA) \\ &\quad + b_{BB}(i+1) P(BB|BB) P(x_{i+1} | BB) \\ &\quad + b_{AB}(i+1) P(AB|BB) P(x_{i+1} | AB) \\ &= p b_{BB}(i+1) P(x_{i+1} | BB) \\ &\quad + (1-p) b_{AB}(i+1) P(x_{i+1} | AB) \end{aligned}$$

$$\begin{aligned}
b_{AB}(i) &= b_{AA}(i+1) P(AA|AB) P(x_{i+1}|AA) \\
&\quad + b_{BB}(i+1) P(BB|AB) P(x_{i+1}|BB) \\
&\quad + b_{AB}(i+1) P(AB|AB) P(x_{i+1}|AB) \\
&= P_A b_{AA}(i+1) P(x_{i+1}|AA) \\
&\quad + P_B b_{BB}(i+1) P(x_{i+1}|BB) \\
&\quad + P b_{AB}(i+1) P(x_{i+1}|AB)
\end{aligned}$$

Initialization

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

Posterior decoding



$$P(z_i | y_1 \dots y_L) = \frac{P(x_1 \dots x_i z_i \dots x_L)}{P(y_1 \dots y_L)}$$

$$= \frac{P(x_1 \dots x_i z_i) P(x_{i+1} \dots x_L | z_i)}{P(y_1 \dots y_L)}$$

$$= \frac{f_{z_i}(i) b_{z_i}(i)}{P(y_1 \dots y_L)}$$

$$P(z_i = AA | y_1 \dots y_L) = \frac{f_{AA}(i) b_{AA}(i)}{f_{AA}(i) b_{AA}(i) + f_{BB}(i) b_{BB}(i) + f_{AB}(i) b_{AB}(i)}$$

$$P(z_i = BB | y_1 \dots y_L) = \frac{f_{BB}(i) b_{BB}(i)}{\text{same}}$$

$$P(z_i = AB | y_1 \dots y_L) = \frac{f_{AB}(i) b_{AB}(i)}{\text{same}}$$