

Genetic Data and EM of Li (2011)

Sequencing Data

ATTG(ATT)GC
ATTG(C)TGC

↑ SNP = location where ^{"allele"} differences occur

Typically, SNPs are "biallelic" (only two alleles occur in population)

Choose one allele as the "reference" (doesn't matter which)

The individual's "dosage" at a SNP is the number of reference alleles that individual has

Let b_{i1} = individual i 's dosage @ SNP 1
 b_{i2} = individual i 's dosage @ SNP 2

$$b = (b_1, b_2), \quad b_i = (b_{i1}, b_{i2})$$

b_1, b_2, \dots, b_n iid b

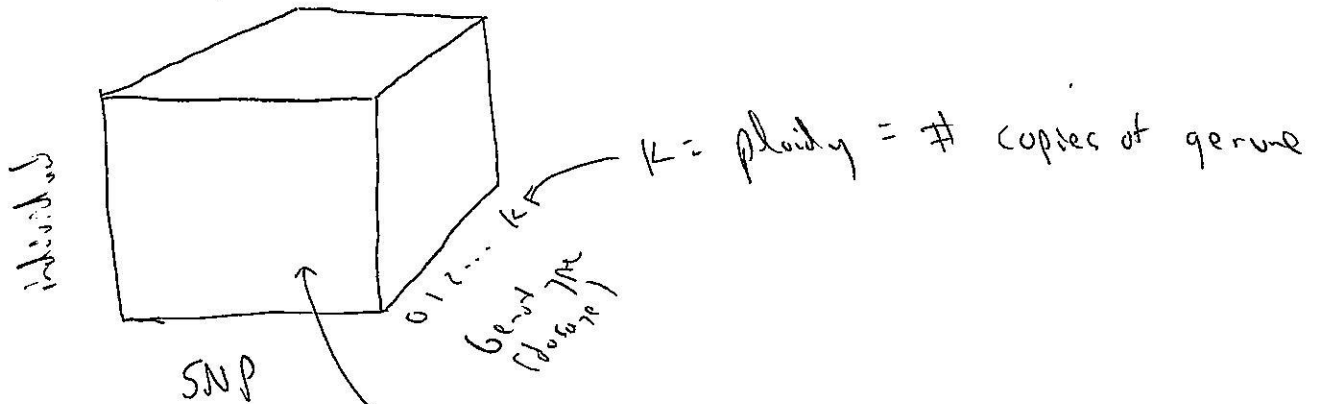
Goal: estimate $\rho = \text{cor}(b_1, b_2)$ using sample b_1, b_2, \dots, b_n

If we observed b_1, b_2, \dots, b_n then this would be easy

Sample Correlation:
$$\hat{\rho} = \frac{\sum_{i=1}^n (b_{i1} - \bar{b}_{1.})(b_{i2} - \bar{b}_{2.})}{\sqrt{\sum_{i=1}^n (b_{i1} - \bar{b}_{1.})^2} \sqrt{\sum_{i=1}^n (b_{i2} - \bar{b}_{2.})^2}}$$

We do not observe $G_{i,j}, G_n$

We observe Genotype likelihoods for each individual at each genotype for each locus



$$\Pr(\text{data} \mid G_{j,i} = k)$$

$j = \text{SNP}$

$i = \text{individual}$

$k = \text{dosage (Genotype)}$

Considering just 2 loci at a time, let

$$a_{i1} = \Pr(\text{Data}_1 \mid G_{i1} = k)$$

$$b_{i2} = \Pr(\text{Data}_2 \mid G_{i2} = k)$$

We sometimes have prior information on distribution of genotypes:

$$\pi_{k1} = \Pr(G_1 = k) \leftarrow \text{Proportion individuals w/ genotype } k \text{ at locus } 1$$

$$\pi_{k2} = \Pr(G_2 = k)$$

Recall Bayes Rule:

$$\Pr(G_{ii} = k | \text{Data}_i) = \frac{\Pr(\text{Data}_i | G_{ii} = k) \Pr(G_{ii} = k)}{\Pr(\text{Data}_i)}$$

$$\Pr(\text{Data}_i | G_{ii} = k) = \text{genotype likelihood} = a_{ik}$$

$$\Pr(G_{ii} = k) = \pi_{ki}$$

$$\begin{aligned} \Pr(\text{Data}_i) &= \sum_{k=0}^K \Pr(\text{Data}_i | G_{ii} = k) \Pr(G_{ii} = k) \\ &= \sum_{k=0}^K a_{ik} \pi_{ki} \end{aligned}$$

So posterior probability of individual i 's genotype given data is

$$\Pr(G_{ii} = k | \text{Data}_i) = \frac{a_{ik} \pi_{ki}}{\sum_{k=0}^K a_{ik} \pi_{ki}}$$

What if we don't know π ?

↑ we can estimate it from the data!

$$P(\text{All data} | \pi) = \prod_{i=1}^n P(\text{Data}_i | \pi) \\ = \prod_{i=1}^n \left(\sum_{k=0}^K a_{ik} \pi_k \right)$$

log

observed function

$$\sum_{i=1}^n \log \left[\sum_{k=0}^K a_{ik} \pi_k \right]$$

↑ Maximum likelihood estimation says to estimate π by maximizing this quantity over π

Hard to do this, so we use EM algorithm
"Expectation - Maximization"

$$P(\text{Data}_i \text{ and } G_i | \pi) = \prod_{k=0}^K [a_{ik} \pi_k]^{\mathbb{1}(G_i=k)}$$

Indicator function

log

$$P(\text{All data and All } G | \pi) = \prod_{i=1}^n \prod_{k=0}^K (a_{ik} \pi_k)^{\mathbb{1}(G_i=k)} \\ \sum_{i=1}^n \sum_{k=0}^K \mathbb{1}(G_i=k) [\log(a_{ik}) + \log(\pi_k)]$$

↑ we don't know G_i , but if we knew π_k then we could know the distribution of G_i | data_i

↑ then we could maximize this "expected log likelihood" to get a new π

EM: start w π_0

$$E\text{-step: } E[1(G:=h) \mid \text{data}, \pi_0] = P(G:=h \mid \text{data}, \pi_0)$$

$$= \frac{a_{ih} \pi_0}{\sum_{k=0}^K a_{ik} \pi_0} =: w_{ih}$$

$$M\text{-step: } \underset{\pi}{\text{maximize}} \sum_{i=1}^n \sum_{h=0}^K w_{ih} [\log(a_{ih}) + \log(\pi_h)]$$

$$= \underset{\pi}{\text{maximize}} \sum_{i=1}^n \sum_{h=0}^K w_{ih} \log(\pi_h)$$

$$= \underset{\pi}{\text{maximize}} \sum_{h=0}^K \log(\pi_h) \underbrace{\sum_{i=1}^n w_{ih}}_{=: w_h}$$

$$\frac{d}{d\pi_h} \left(\sum_{h=0}^K \log(\pi_h) w_h + \lambda \left(\sum_{h=0}^K \pi_h - 1 \right) \right)$$

$$\frac{w_h}{\pi_h} + \lambda \stackrel{!}{=} 0$$

$$\Rightarrow \pi_h = -\frac{w_h}{\lambda}$$

$\Rightarrow \pi_h \propto w_h$ and must sum to 1

$$\Rightarrow \pi_h = \frac{w_h}{\sum w_h}$$