# FFR110 - Homework Problem 3 - Part 2

Computational Biology

**Felix Waldschock**
000420-T398

Chalmers University of Technology
February 2024

# Contents

# 1 Population genetics

**1.1** **a) Derive an expression for P(Sn = 0), i.e. the probability to not have any SNPs in a sample of size n.**

# Problem 2 a

$F_2$ describes 2 alleles to

$$F_2^{(t+1)} = (1-\mu)^2 \cdot \left[ \frac{1}{N} + \left(1 - \frac{1}{N}\right) F_2^{(t)} \right]$$

in prev. generation

prob
2 child
are same

no mutation

Same
parent

diffrent
parents

parents
are same

for $\mu = 0$ and $N \to \infty$

this leads to:

$$\overline{F_2} = \frac{1}{1+\theta}$$

What we want to find is $P(S_n = 0)$ which describes
the probability of $n$ alleles not being the same.

from the lecture notes:

$$P(S=j) = \frac{(\mu T_c)^j}{j!} e^{-\mu T_c} \qquad \text{where } T_c = \sum_{j=2}^{n} j\, T_j \qquad (1)$$

number of
possible pairs

$$\text{and } P(T_j) = \lambda_j e^{-\lambda_j T_j} \qquad \text{with } \lambda_j = \frac{\binom{j}{2}}{N} \qquad (2)$$

Time to coalescent event backward
in time when starting with $j$ lines.

now  $P(S_n = 0)$ :

$$P(S_n = 0) = e^{\mu T_c} = e^{\mu \sum_{j=2}^{n} j \cdot T_j}$$

rewrite as product:

$$\prod_{n \geq 2} e^{-\mu T_n n}$$

now integrate this probability from $0 \to \infty$ to get $F_n$.

$$F_n = \prod_{n \geq 2} \int_0^\infty e^{-\mu \cdot T_n \cdot n} \cdot P(T_n) \cdot dT_n$$

now insert (2)

$$= \prod_{n \geq 2} \int_0^\infty e^{-\mu T_n n} \cdot \frac{1}{N} \binom{n}{2} e^{-\binom{n}{2}\frac{T_n}{N}} \, dT_n$$

$$= \prod_{n \geq 2} \binom{n}{2} \int_0^\infty e^{-\mu T_n n} \cdot e^{-\binom{n}{2}\frac{T_n}{N}} \frac{1}{N} \, dT_n$$

simplify the exp expression with $T_n = N \cdot t$ ; $e^{-\mu N t n}$ & $e^{-\binom{n}{2} t}$ ; $dT_n \to dt$

$$= \prod_{n \geq 2} \binom{n}{2} \int_0^\infty e^{-t\left(\mu N n + \binom{n}{2}\right)} \, dt \; .$$

now:

$$\binom{n}{2} = \frac{n(n-1)}{2} \; ; \quad \prod_{n \geq 2} \binom{n}{2} = n-1 \qquad\qquad (3)$$

solve the integral: $\int_0^\infty e^{-nx} = -\dfrac{e^{-nx}}{n}$ gives:

$$\prod_{n \geqslant 2} \left( -\frac{\binom{n}{2}}{\mu N n + \binom{n}{2}} \left[ e^{-\left(\mu N n + \binom{n}{2}\right)t} \right]_0^{-\infty} \right)$$

with (3) this gives:

$$= \prod_{n \geqslant 2} \frac{n(n-1)\cancel{(n-2)!}}{\cancel{n}\theta\cancel{(n-2)!} + \cancel{n}(n-1)\cancel{(n-2)!}}$$

$$= \prod_{n \geqslant 2} \frac{(n-1)}{\theta + n - 1}$$

$$\prod \frac{n-1}{\theta + n - 1} = \frac{(n-1)!}{(1+\theta)(2+\theta)\ldots(n-1+\theta)}$$

## 1.2 b) Derive the distribution of the number of SNPs in a sample of size n = 2.

$$P(S_2 = j) = \frac{1}{1+\theta} \left( \frac{\theta}{1+\theta} \right)^j \tag{1}$$

# Problem 2b

Derive distribution of the number of SNPs in a sample of size $n=2$.

should find:

$$P(S_2 = j) = \frac{1}{1+\theta} \left(\frac{\theta}{1+\theta}\right)^j \qquad (2b)$$

from lecture notes:

$$T_c = \sum_{j=2}^{n} j\, T_j \quad ; \quad \text{with} \quad S_2 = j \implies T_c = 2 T_2$$

$$P(T_j) = \lambda_j \cdot e^{-\lambda_j T_j} \quad \text{with} \quad \lambda_j = \frac{\binom{j}{2}}{N} \quad ; \quad j=2 \implies P(T_j) = \frac{1}{N} \cdot e^{-T_2}$$

$F_2$ is found when integrating the product of the 2 above:

$$F_2 = \int_0^\infty \underbrace{P(S_2 = j)}_{\substack{\text{expression} \\ \text{from a)}}} \cdot \underbrace{P(T_2)}_{\substack{\text{also} \\ \text{see a)}}} dT_2 = \int_0^\infty \frac{(2N\,T_2)^j}{j!}\, e^{-\frac{(2\mu N + 1) T_2}{N}} \cdot \frac{1}{N} \cdot dT_2$$

transform $T_2 = N \cdot t$ leads to:

$$= \int_0^\infty \frac{(2\mu N \cdot t)^j}{j!} \cdot e^{-(2\mu N + 1) \cdot t} \cdot dt$$

now we rewrite $2\mu T_2$

$$= \int_0^\infty \frac{(\theta t)^j}{j!}\, e^{-(\theta + 1) t}\, dt$$

$$= \frac{\theta^j}{j!} \int_0^\infty t^j\, e^{-(\theta + 1) t}\, dt$$

$\theta = 2 N \cdot \mu$     (lecture notes p. 7)

‖ take $t$ independent variables out of $\int$

$$= \frac{\theta^j}{j!} \left( \left[ -\frac{(e^{-(\theta+\lambda)t} \cdot t^j)}{\theta+\lambda} \right]_0^\infty + \frac{j}{\theta+j} \int_0^\infty t^{j-1} e^{-(\theta+\lambda)t} dt \right)$$

this will give new integrals until $j=0$ aut $t^j = 1$

$$= -\frac{\theta^j}{j! \, (\theta+\lambda)} \left[ \sum_{k=1}^{j} e^{-(\theta+\lambda)t} t^k \right]_0^\infty + \left( \frac{\theta}{\theta+\lambda} \right)^j \int_0^\infty e^{-(\theta+\lambda)t} dt$$

$$\underbrace{\hspace{5cm}}_{(\#)} \qquad \underbrace{\hspace{3cm}}_{(*)}$$

as we want to find (2b) there are some terms that
seem familiar. Therefore compute $(*)$.

$$\int_0^\infty e^{-(\theta+\lambda)t} dt = -\frac{1}{\theta+\lambda} \left[ e^{-(\theta+\lambda)t} \right]_0^\infty = \frac{1}{1+\theta}$$

$$\underbrace{\hspace{3cm}}_{1/0}$$

this looks again familiar.

Therefore lets see if $(\#)$ cancels out.

$$-\frac{\theta^j}{j! \, (\theta+\lambda)} \left[ \sum_{k=1}^{j} e^{-(\theta+\lambda)t} t^k \right]_0^\infty$$

$$\underbrace{\hspace{2cm}}_{!0} \qquad \underbrace{\hspace{3cm}}_{?0}$$

$$\left[ \sum_{k=1}^{j} e^{-(\theta+\lambda)t} t^k \right]_0^\infty = \lim_{t \to \infty} \sum_{k=1}^{j} e^{-(\theta+\lambda)t} t^j - \lim_{t \to 0} \sum_{k=1}^{j} e^{-(\theta+\lambda)t} t^k$$

$$\underbrace{\hspace{4cm}}_{=0}$$

$$\lim_{t \to \infty} \sum_{k=1}^{j} e^{-(\theta+\lambda)t} t^k = \sum_{k=1}^{j} \lim_{t \to \infty} e^{-(\theta+\lambda)t} t^k$$

new look with l'Hospital

$$t^k \xrightarrow{\;'\;} k \cdot t^{k-1}$$

$$e^{-(\theta+\lambda)t} = \frac{1}{e^{-(\theta+\lambda)t}} \xrightarrow{\;'\;} \left((\theta+\lambda)\left(e^{(\theta+\lambda)t}\right)\right)^{-\lambda}$$

$$\sum_{k=1}^{j} \lim_{t \to \infty} \frac{k \cdot t^{k-1}}{(\theta+\lambda)\left(e^{(\theta+\lambda)t}\right)}$$

and again l'Hospital ... leads to

$$\sum_{k=1}^{j} \lim_{t \to \infty} \frac{i!}{(\theta+\lambda)^i \left(e^{(\theta+\lambda)t}\right)} = \underline{0}$$

So this then shows:

$$= -\frac{\theta^j}{j! \,(\theta+\lambda)} \underbrace{\left[\sum_{k=1}^{j} e^{-(\theta+\lambda)t} \; t^k\right]_0^{\infty}}_{=0} + \left(\frac{\theta}{\theta+\lambda}\right)^j \int_0^{\infty} e^{-(\theta+\lambda)t}\, dt$$

$$= \underline{\underline{\left(\frac{\theta}{\theta+j}\right)^j \cdot \frac{1}{1+\theta}}}$$