

## Revisiting allelic frequencies estimation: A decision theory approach to derive Bayes, minimax and admissible estimators

Carlos Alberto Martínez Niño<sup>1,2</sup>, Kshitij Khare<sup>2</sup>,  
Mauricio A. Elzo<sup>1</sup>

Department of Animal Sciences<sup>1</sup>  
Department of Statistics<sup>2</sup>  
University of Florida



## Problem and objective

- Necessity of deriving point estimators of allelic frequencies with appealing statistical properties and biological soundness.
- Are allelic frequencies unknown constants or random variables?
- Random variation of allelic frequencies due to certain evolutionary forces (Wright, 1930; 1937).
- The aim of this study was to derive alternative estimators of allele frequencies with optimal statistical properties under a decision theory framework.



Vol. 30, 1937 GENETICS: S. WRIGHT  
THE DISTRIBUTION OF GENE FREQUENCIES IN POPULATIONS  
By Sewall Wright

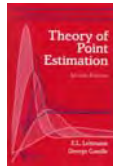
## Elements of decision theory

- Parameter space  $\theta$ , decision space  $D$ , observed data  $\mathbf{X}$ , loss function  $L(\theta, \delta(\mathbf{X}))$
- Frequentist risk:  $R(\theta, \delta) = E_{\theta}[L(\theta, \delta(\mathbf{X}))]$ .
- Decision rules with uniformly smallest risk rarely exist (Lehmann and Casella, 1998): Use a weaker optimality criterion.
- Bayes decision rule:

$$r(\mathbf{A}, \delta^*) = \int_{\theta} R(\theta, \delta^*) d\Lambda(\theta) = \inf_{\delta \in D} r(\mathbf{A}, \delta).$$

- Minimax decision rule:

$$\sup_{\theta \in \Theta} R(\theta, \delta^*) = \inf_{\delta \in D} \sup_{\theta \in \Theta} R(\theta, \delta)$$

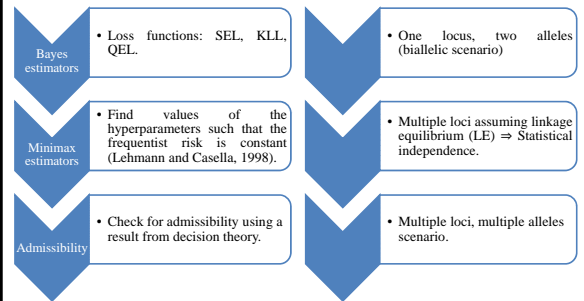


$$P(\mathbf{A}|\theta) = \frac{P(\theta|\mathbf{A})P(\mathbf{A})}{P(\theta)}$$

## Approach

Statistical level

Statistical-genetics level



## General setting and notation

Let  $X_1, X_2$  and  $X_3$  be random variables indicating the number of animals having genotypes AA, AB and BB and assume Hardy-Weinberg equilibrium.

Let  $\theta :=$  frequency of the “reference” allele B.

$$\mathbf{X} := (X_1, X_2, X_3)$$

$$\mathbf{X}|\theta \sim \text{Trinomial}(n; (1-\theta)^2, 2\theta(1-\theta), \theta^2)$$

$$\theta \sim \text{Beta}(\alpha, \beta)$$

## Bayes estimators and risks

Loss function	Functional form of loss function	Bayes estimator (BE)	Functional form of BE	Frequentist risk
SEL	$(\theta - \delta)^2$	Posterior mean	$\frac{x_2 + 2x_3 + \alpha}{2n + \alpha + \beta}$	$\frac{2n\theta(1-\theta) + [\alpha(1-\theta) - \beta\theta]^2}{(2n + \alpha + \beta)^2}$
KLL	$E_{\theta} \left[ \ln \left( \frac{\pi(\mathbf{X} \theta)}{\pi(\mathbf{X} \delta)} \right) \right]$	$\text{argmin}_{\delta \in D} \int_0^1 L_{KLL}(\theta, \delta) \pi(\theta \mathbf{X}) d\theta$	$\frac{x_2 + 2x_3 + \alpha}{2n + \alpha + \beta}$	No closed form
QEL	$\frac{(\theta - \delta)^2}{\theta(1-\theta)}$	Mean of: $w(\theta)\pi(\theta \mathbf{X})$	See next slide	See next slide

## Bayes estimators and risks: QEL

$$\begin{aligned} \int_0^1 w(\theta) (\theta - \hat{\theta}^{QEL})^2 \pi(\theta|X) d\theta &\propto \int_0^1 (\theta - \hat{\theta}^{QEL})^2 \theta^{\alpha-2} (1-\theta)^{2x_1+\beta-2} d\theta \\ &= \int_0^1 \theta^{\alpha} (1-\theta)^{2x_1+\beta-2} d\theta - 2\hat{\theta}^{QEL} \int_0^1 \theta^{\alpha-1} (1-\theta)^{2x_1+\beta-2} d\theta \\ &\quad + (\hat{\theta}^{QEL})^2 \int_0^1 \theta^{\alpha-2} (1-\theta)^{2x_1+\beta-2} d\theta \end{aligned}$$

Finite iff  $\hat{\theta}^{QEL} = 0$ .

$$\hat{\theta}^{QEL} = \begin{cases} \frac{x_2+2x_3+\alpha-1}{2n+\alpha+\beta-2}, & \text{if } x_2+2x_3+\alpha-1 > 0, 2x_1+x_2+\beta-1 > 0 \\ 0, & \text{if } x_2+2x_3+\alpha-1 \leq 0 \\ 1, & \text{if } 2x_1+x_2+\beta-1 \leq 0 \end{cases}$$

$$R(\theta, \hat{\theta}^{QEL}) = \begin{cases} \frac{2n}{(2n+\alpha+\beta-2)^2} + \frac{(-\theta(\alpha+\beta-2)+\alpha-1)^2}{\theta(1-\theta)(2n+\alpha+\beta-2)^2}, & \text{if } x_2+2x_3+\alpha-1 > 0, 2x_1+x_2+\beta-1 > 0 \\ \frac{\theta}{1-\theta}, & \text{if } x_2+2x_3+\alpha-1 \leq 0 \\ \frac{1-\theta}{\theta}, & \text{if } 2x_1+x_2+\beta-1 \leq 0 \end{cases}$$

## Derivation of minimax rules

*Theorem 1* (Lehmann and Casella, 1998). Let  $\Lambda$  be a prior and  $\delta_\Lambda$  a Bayes rule with respect to  $\Lambda$  with Bayes risk satisfying  $r(\Lambda, \delta_\Lambda) = \sup_{\theta \in \Theta} R(\theta, \delta_\Lambda)$ . Then: i)  $\delta_\Lambda$  is minimax and ii)  $\Lambda$  is least favorable.

Loss function	Hyperparameters	Functional form of BE	Frequentist risk
SEL	$\alpha = \frac{\sqrt{n}}{2}, \beta = \frac{\sqrt{n}}{2}$	$\frac{x_2+2x_3+\frac{\sqrt{n}}{2}}{\sqrt{2n}(\sqrt{2n}+1)}$	$(4(1+\sqrt{2n})^2)^{-1}$
QEL <sup>1</sup>	$\alpha = 1, \beta = 1$	$\frac{x_2+2x_3}{2n} = MLE$	$\frac{1}{2n}$

<sup>1</sup> Provided  $x_2+2x_3+\alpha-1 > 0, 2x_1+x_2+\beta-1 > 0$ .

## Extension to k loci (LE and independent priors)

$\theta = (\theta_1, \theta_2, \dots, \theta_k)$ .  $X = (X_1, X_2, \dots, X_k)$ .  $X_i = (X_{1i}, X_{2i}, X_{3i})$

Minimize:  $R(\theta, \delta) = \int_{\theta_1} \dots \int_{\theta_k} L(\theta, \delta(X)) \pi(\theta|X) d\theta_1 \dots d\theta_k$ , wrt  $\delta_i, \forall i = 1, 2, \dots, k$

$$\begin{aligned} R(\theta, \delta) &= \int_{\theta_1} \dots \int_{\theta_k} \left( \sum_{i=1}^k L(\theta_i, \delta_i(X)) \right) \pi(\theta|X) d\theta_1 \dots d\theta_k \\ &= \sum_{i=1}^k \int_{\theta_1} \dots \int_{\theta_k} L(\theta_i, \delta_i(X)) \prod_{j=1}^k \pi(\theta_j|X_j) d\theta_1 \dots d\theta_k \end{aligned}$$

The  $h^{\text{th}}$  integral in the summation ( $h = 1, 2, \dots, k$ ) can be written as:

$$\begin{aligned} \int_{\theta_h} L(\theta_h, \delta_h(X)) \pi(\theta_h|X_h) d\theta_h &\times \int_{\theta_1} \dots \int_{\theta_{h-1}} \int_{\theta_{h+1}} \dots \int_{\theta_k} \prod_{j \neq h} \pi(\theta_j|X_j) d\theta_1 \dots d\theta_{h-1} d\theta_{h+1} \dots d\theta_k \\ &= \int_{\theta_h} L(\theta_h, \delta_h) \pi(\theta_h|X_h) d\theta_h \end{aligned}$$

Bayes estimation of  $\theta$  reduces to that of its components.

## Multiallelic loci

Let  $\theta_{1i}, \theta_{2i}, \dots, \theta_{n_i}$  be the frequencies of the  $n_i$  alleles of locus  $i$ .

Let  $Y_{1i}, Y_{2i}, \dots, Y_{n_i}$  be random variables indicating the counts of each one of the  $n_i$  allelic variants at locus  $i, i = 1, 2, \dots, k$ .

$Y_i := (Y_{1i}, Y_{2i}, \dots, Y_{n_i}) \sim \text{Multinomial}(2n, \theta_i)$

$\theta_i = (\theta_{1i}, \theta_{2i}, \dots, \theta_{n_i}) \sim \text{Dirichlet}(\alpha_i = (\alpha_{1i}, \alpha_{2i}, \dots, \alpha_{n_i}))$

$\therefore \theta_i | y_i \sim \text{Dirichlet}(\alpha_{1i} + y_{1i}, \alpha_{2i} + y_{2i}, \dots, \alpha_{n_i} + y_{n_i})$

Under the loss  $\sum_{j=1}^{n_i} (\hat{\theta}_{ji} - \theta_{ji})^2$ ,  $\hat{\theta}_i^{M-SEL} = (\hat{\theta}_{ji})_{n_i \times 1} = \frac{\alpha_{ji} + y_{ji}}{2n + \sum_{j=1}^{n_i} \alpha_{ji}}$

Under the loss  $\sum_{j=1}^{n_i} \theta_{ji}^{-1} (\hat{\theta}_{ji} - \theta_{ji})^2$ :

$$\hat{\theta}_i^{M-QEL} = (\hat{\theta}_{ji}^{M-QEL})_{n_i \times 1} = \begin{cases} \frac{\alpha_{ji} + y_{ji} - 1}{\sum_{j=1}^{n_i} \alpha_{ji} + 2n - 1}, & \text{if } \alpha_{ji} + y_{ji} - 1 > 0 \\ 0, & \text{if } \alpha_{ji} + y_{ji} - 1 \leq 0 \end{cases}$$

## Admissibility

Admissibility of one-dimensional and vector-valued estimators was established using this theorem (Lehmann and Casella, 1998).

*Theorem 2.* For a possibly vector-valued parameter  $\theta$ , suppose that  $\delta^\pi$  is a Bayes estimator having finite Bayes risk with respect to a prior density  $\pi$  which is positive for all  $\theta \in \Theta$ , and that the risk function of every estimator  $\delta$  is a continuous function of  $\theta$ . Then  $\delta^\pi$  is admissible.

## Results and comments

$\hat{\theta}^{SEL}$ ,  $\hat{\theta}^{Minimax_1}$  and  $\hat{\theta}^{Minimax_2}$  are admissible, for  $\hat{\theta}^{QEL}$  the property holds provided  $\alpha > 1, \beta > 1$ .

If both alleles are observed: the MLE is also minimax and admissible. We have a Bayes, minimax, admissible and unbiased estimator.

$\hat{\theta}^{M-SEL}$ ,  $\hat{\theta}^{M-Minimax_1}$  and  $\hat{\theta}^{M-Minimax_2}$  are admissible, as well as  $\hat{\theta}^{M-QEL}$  when  $\alpha_{ji} > 1, \forall j_i = 1, 2, \dots, n_i, \forall i = 1, 2, \dots, k$ .

The estimators proposed here always have uniformly smaller variance than the MLE, except for those derived from QEL which require:  $\alpha + \beta > 2$  (biallelic case) and  $\sum_{k_i=1}^{n_i} \alpha_{k_i} > 1$  (multiallelic case) to meet this property.

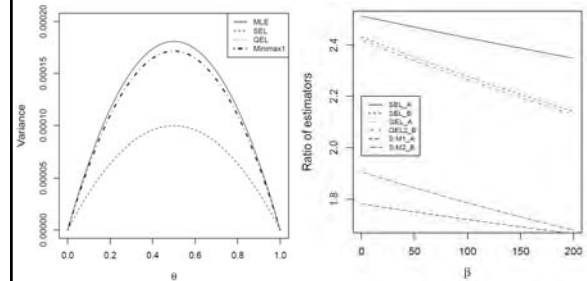
## Frequentist variances

Biallelic case	Multiallelic case
$Var_{\theta}[\hat{\theta}^{ML}] = \frac{\theta(1-\theta)}{2n}$	$Var_{\theta_{j_i}}[(\hat{\theta}_i^{ML})_j] = \frac{\theta_{j_i}(1-\theta_{j_i})}{2n}$
$Var_{\theta}[\hat{\theta}^{SEL}] = \frac{2n\theta(1-\theta)}{(2n+\alpha+\beta)^2}$	$Var_{\theta_{j_i}}[(\hat{\theta}_i^{M-SEL})_j] = \frac{2n\theta_{j_i}(1-\theta_{j_i})}{(2n+\alpha')^2}$
$Var_{\theta}[\hat{\theta}^{Mimax_1}] = \frac{\theta(1-\theta)}{(\sqrt{2n}+1)^2}$	$Var_{\theta_{j_i}}[(\hat{\theta}_i^{M-Mimax_1})_j] = \frac{\theta_{j_i}(1-\theta_{j_i})}{(\sqrt{2n}+1)^2}$
$Var_{\theta}[\hat{\theta}^{Mimax_2}] = Var_{\theta}[\hat{\theta}^{ML}]$	$Var_{\theta_{j_i}}[(\hat{\theta}_i^{M-Mimax_2})_j] = \frac{2n\theta_{j_i}(1-\theta_{j_i})}{(2n+n_i-1)^2}$
If $x_2+2x_3+\alpha-1 > 0$ , $2x_1+x_2+\beta-1 > 0$ : $Var_{\theta}[\hat{\theta}^{QEL}] = \frac{2n\theta(1-\theta)}{(2n+\alpha+\beta-2)^2}$	If $\alpha_{j_i}+y_{j_i}-1 > 0$ : $Var_{\theta_{j_i}}[(\hat{\theta}_i^{M-QEL})_j] = \frac{2n\theta_{j_i}(1-\theta_{j_i})}{(2n+\alpha'-1)^2}$

## Example

Beta distribution with  
 $\alpha = 240, \beta = 240, n = 691$

Sample sizes: 1382 (A), 691 (B)  
 $(x_2, x_3) = (10, 25), \alpha = 96$



## Results and comments

- For all decision rules derived from SEL, the form of the risk functions shows that they converge to zero as  $n \rightarrow \infty$ . QEL: When all hyperparameters are greater than one, all the derived risk functions converge to zero as  $n \rightarrow \infty$ . When some alleles are not observed and the hyperparameters corresponding to their frequencies are smaller or equal to one, the result does not hold.
- The impact of the use of these estimators in the many applications they could have should be assessed either empirically or theoretically and is an area for further research.

## Acknowledgements

- Dr. Malay Ghosh, Department of Statistics, University of Florida.
- Fulbright Colombia and COLCIENCIAS.

