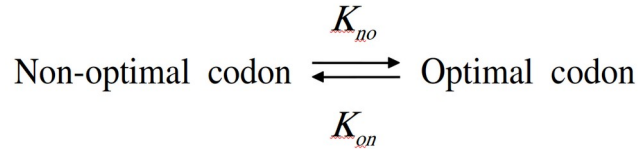


## Estimating the strength of selection on synonymous codon usage

The frequency of optimal codons ( $F_{op}$ ) reflects the balance between the optimal to non-optimal codons synonymous substitution rate ( $K_{on}$ ) and the non-optimal to optimal codons synonymous substitution rate ( $K_{no}$ ):



Substitution rates depend on the corresponding mutation rates ( $\mu_{no}, \mu_{on}$ ) and fixation probabilities ( $P_{no}, P_{on}$ ):

$$K_{no} = 2 N_e \mu_{no} P_{no}$$

$$K_{on} = 2 N_e \mu_{on} P_{on}$$

where  $N_e$  is the effective population size.

Fixation probabilities are given by:

$$P_{no} = \frac{1 - e^{-4 N_e f_0 s}}{1 - e^{-4 N_e s}} = \frac{1 - e^{-2s}}{1 - e^{-4 N_e s}} \xrightarrow{s \rightarrow 0} \frac{2s}{1 - e^{-4 N_e s}} \quad \text{and} \quad P_{on} \xrightarrow{s \rightarrow 0} \frac{-2s}{1 - e^{4 N_e s}}$$

Where  $s$  is the selection coefficient in favor of optimal codons and  $f_0$  the allele frequency of a new arrival mutation ( $f_0 = 1/2N_e$ ).

At equilibrium, the frequency of optimal codons is given by:

$$F_{op} = \frac{K_{no}}{K_{no} + K_{on}}$$

which can be written as:

$$F_{op} = \frac{2 N_e \mu_{no} P_{no}}{2 N_e \mu_{no} P_{no} + 2 N_e \mu_{on} P_{on}}$$

$$F_{op} = \frac{\mu_{no} P_{no}}{\mu_{no} P_{no} + \mu_{on} P_{on}} = \frac{\frac{\mu_{no}}{\mu_{on}} \frac{2s}{1 - e^{-4 N_e s}}}{\frac{\mu_{no}}{\mu_{on}} \frac{2s}{1 - e^{-4 N_e s}} + \frac{-2s}{1 - e^{4 N_e s}}}$$

Let us note lambda, the ratio of mutation rates:  $\lambda = \frac{\mu_{no}}{\mu_{on}}$

$$Fop = \frac{\lambda}{\lambda + \frac{-(1 - e^{-4N_e s})}{1 - e^{4N_e s}}}$$

$$\frac{1}{Fop} = 1 + \frac{1}{\lambda} \times \frac{-(1 - e^{-4N_e s})}{1 - e^{4N_e s}}$$

$$\frac{1}{Fop} - 1 = \frac{1}{\lambda} \times \frac{-(1 - e^{-4N_e s})}{1 - e^{4N_e s}}$$

$$\frac{1 - Fop}{Fop} \times \lambda = \frac{-(1 - e^{-4N_e s})}{1 - e^{4N_e s}}$$

$$\frac{Fop}{1 - Fop} \times \frac{1}{\lambda} = \frac{1 - e^{4N_e s}}{-(1 - e^{-4N_e s})} \quad (1)$$

With the following simplification:

$$\frac{1 - e^{4N_e s}}{-(1 - e^{-4N_e s})} = \frac{1 - e^{4N_e s}}{-\left(1 - \frac{1}{e^{4N_e s}}\right)} = \frac{e^{4N_e s} \times (1 - e^{4N_e s})}{(1 - e^{4N_e s})} = e^{4N_e s}$$

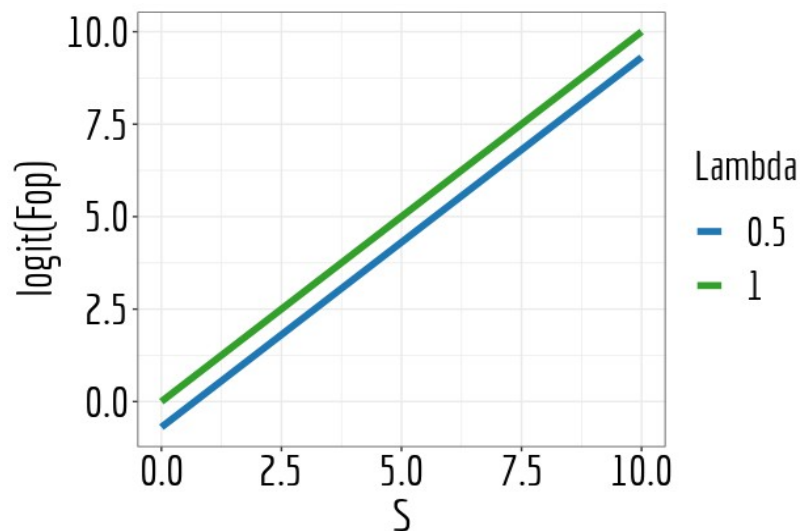
$$(1) \rightarrow \frac{Fop}{1 - Fop} \times \frac{1}{\lambda} = e^{4N_e s}$$

Thus, the population-scaled selection coefficient ( $S = 4N_e s$ ) is given by:

$$S = \log\left(\frac{Fop}{1 - Fop}\right) - \log(\lambda) = \text{logit}(Fop) - \log(\lambda)$$

Hence, we expect a linear correlation between  $\text{logit}(Fop)$  and  $S$ :

$$\text{logit}(Fop) = S + \log(\lambda)$$



If for weakly-expressed genes there is no selection, implied by the non-variation of  $Fop$  with gene expression,  $S_{low-exp} \sim 0$  :

$$\text{logit}(Fop_{low-exp}) = 0 + \log(\lambda)$$

$$\text{logit}(Fop_{high-exp}) = S_{high-exp} + \log(\lambda)$$

$$S_{high-exp} = \text{logit}(Fop_{high-exp}) - \text{logit}(Fop_{low-exp})$$

