

Supplementary Material

1. Derivation of the Distribution of Φ Estimates

Using least squares or maximum likelihood methods, we infer that the estimates of the kinetic parameters follow (asymptotically) a multivariate normal distribution. As the chevron curves for the two amino acid sequences are fit separately, kinetic parameter estimates from the same chevron fit will be correlated, while kinetic parameter estimates from different chevron fits will be independent. Mathematically, this means that

$$\mathbf{Y} = \begin{bmatrix} \ln(\widehat{k_f}) \\ \ln(\widehat{k_u}) \\ \ln(\widehat{k'_f}) \\ \ln(\widehat{k'_u}) \end{bmatrix} \sim N \left(\begin{bmatrix} \ln(k_f) \\ \ln(k_u) \\ \ln(k'_f) \\ \ln(k'_u) \end{bmatrix}, \begin{bmatrix} (\sigma_f)^2 & \rho\sigma_f\sigma_u & 0 & 0 \\ \rho\sigma_f\sigma_u & (\sigma_u)^2 & 0 & 0 \\ 0 & 0 & (\sigma'_f)^2 & \rho'\sigma'_f\sigma'_u \\ 0 & 0 & \rho'\sigma'_f\sigma'_u & (\sigma'_u)^2 \end{bmatrix} \right) \quad (1)$$

Here, $\ln(k_f)$ denotes the logarithm of the folding rate of the wild type, $\ln(k_u)$ denotes the logarithm of the unfolding rate of the wild type, and $\ln(k'_f)$ and $\ln(k'_u)$ denote the respective parameters for the mutant. The symbol σ is used for the respective standard errors, and ρ is used for the respective correlations. Figure 1 shows scatterplots for the kinetic parameter estimates of a chevron fit obtained in a statistical simulation study.

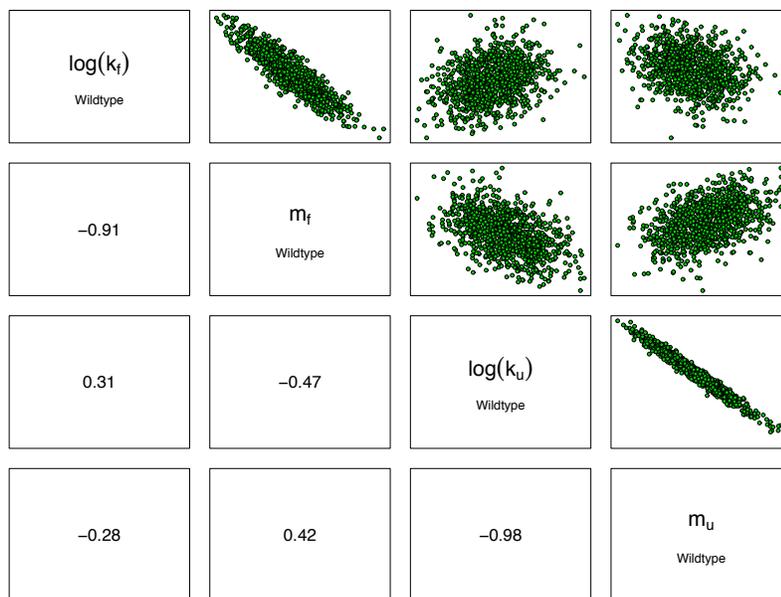


Figure 1: For 10,000 iterations, Gaussian noise with standard deviation typically seen for the experimental error in such kinetic studies (see for example www.foldomics.org) was added to a pair of “synthetic” chevron curves, and the kinetic parameters were estimated by fitting the chevrons. For brevity, the distribution of the kinetic parameter estimates (including the estimates for the folding and unfolding arms, m_f and m_u respectively) are shown for only one chevron. This figure is only intended to highlight the distributions between the kinetic parameter estimates, and thus the axis labelings are omitted.

An elementary result of linear model theory (e. g. Seber [1]) states that if \mathbf{Y} is a random vector of length n following a multivariate normal distribution, i. e. $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and \mathbf{C} is a matrix with n columns and p

rows, then $\mathbf{CY} \sim N_p(\mathbf{C}\boldsymbol{\mu}, \mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^T)$. As the estimates in the changes in free energy are linear functions of the kinetic parameters ($\widehat{\Delta\Delta G}_\ddagger = RT \times [\ln(\widehat{k}_f) - \ln(\widehat{k}'_f)]$, and $\widehat{\Delta\Delta G}_u = RT \times [\ln(\widehat{k}_f) - \ln(\widehat{k}'_f) - \ln(\widehat{k}_u) + \ln(\widehat{k}'_u)]$), we choose

$$\mathbf{C} = RT \times \begin{bmatrix} +1 & 0 & -1 & 0 \\ +1 & -1 & -1 & +1 \end{bmatrix}, \quad (2)$$

where R is the gas constant, and T is the absolute temperature. Therefore

$$\mathbf{CY} = RT \times \begin{bmatrix} +1 & 0 & -1 & 0 \\ +1 & -1 & -1 & +1 \end{bmatrix} \times \begin{bmatrix} \ln(\widehat{k}_f) \\ \ln(\widehat{k}_u) \\ \ln(\widehat{k}'_f) \\ \ln(\widehat{k}'_u) \end{bmatrix} = RT \times \begin{bmatrix} \ln(\widehat{k}_f) - \ln(\widehat{k}'_f) \\ \ln(\widehat{k}_f) - \ln(\widehat{k}'_f) - \ln(\widehat{k}_u) + \ln(\widehat{k}'_u) \end{bmatrix} = \begin{bmatrix} \widehat{\Delta\Delta G}_\ddagger \\ \widehat{\Delta\Delta G}_u \end{bmatrix}, \quad (3)$$

and its distribution is given by

$$\begin{bmatrix} \widehat{\Delta\Delta G}_\ddagger \\ \widehat{\Delta\Delta G}_u \end{bmatrix} \sim N \left(\begin{bmatrix} \Delta\Delta G_\ddagger \\ \Delta\Delta G_u \end{bmatrix}, \begin{bmatrix} \sigma_\ddagger^2 & \rho_\Delta \sigma_\ddagger \sigma_u \\ \rho_\Delta \sigma_\ddagger \sigma_u & \sigma_u^2 \end{bmatrix} \right), \quad (4)$$

where

$$\begin{aligned} \sigma_\ddagger^2 &= \text{var}(\widehat{\Delta\Delta G}_\ddagger) \\ &= (RT)^2 \times [(\sigma_f)^2 + (\sigma'_f)^2], \\ \sigma_u^2 &= \text{var}(\widehat{\Delta\Delta G}_u) \\ &= (RT)^2 \times [(\sigma_f)^2 + (\sigma'_f)^2 + (\sigma_u)^2 + (\sigma'_u)^2 - 2\rho\sigma_f\sigma_u - 2\rho\sigma'_f\sigma'_u], \\ \rho_\Delta \sigma_\ddagger \sigma_u &= \text{cov}(\widehat{\Delta\Delta G}_\ddagger, \widehat{\Delta\Delta G}_u) \\ &= \text{corr}(\widehat{\Delta\Delta G}_\ddagger, \widehat{\Delta\Delta G}_u) \times \text{se}(\widehat{\Delta\Delta G}_\ddagger) \times \text{se}(\widehat{\Delta\Delta G}_u) \\ &= (RT)^2 \times [(\sigma_f)^2 + (\sigma'_f)^2 - \rho\sigma_f\sigma_u - \rho\sigma'_f\sigma'_u]. \end{aligned} \quad (5)$$

Here, $\text{var}()$ denotes the variance, $\text{cov}()$ the covariance, $\text{corr}()$ the correlation, and $\text{se}()$ the standard error of the respective arguments. Figure 2 shows the marginal and joint distributions of $\widehat{\Delta\Delta G}_\ddagger$ and $\widehat{\Delta\Delta G}_u$ derived from our simulation study.

As the Φ -value estimate is not a linear function of the estimates of the changes in free energy, we use a Taylor-series expansion to derive an approximate distribution. This method is also often referred to as *error propagation* in the life sciences, and *Delta method* in the statistical literature. In brief, it states that if we have an estimator $\hat{\boldsymbol{\theta}}$ for a parameter $\boldsymbol{\theta}$, following a multivariate normal distribution with mean $\boldsymbol{\theta}$ and variance covariance matrix $\boldsymbol{\Sigma}$, then for a non-constant and differentiable function f ,

$$f(\hat{\boldsymbol{\theta}}) \approx N(f(\boldsymbol{\theta}), \mathbf{V}). \quad (6)$$

where

$$\mathbf{V} = \left(\frac{\delta f}{\delta \boldsymbol{\theta}} \right)^T \boldsymbol{\Sigma} \left(\frac{\delta f}{\delta \boldsymbol{\theta}} \right) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}. \quad (7)$$

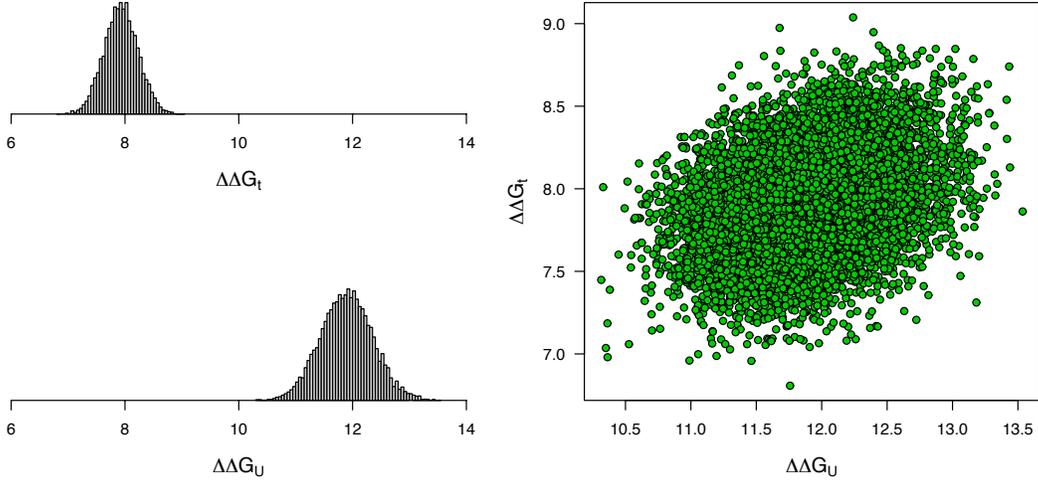


Figure 2: The marginal and joint distributions of $\widehat{\Delta\Delta G_{\ddagger}}$ and $\widehat{\Delta\Delta G_U}$ derived from our simulation study, confirming the theoretical result that $\widehat{\Delta\Delta G_{\ddagger}}$ and $\widehat{\Delta\Delta G_U}$ follow a bivariate normal distribution.

The theorem is often stated as a limit theorem, i. e. in terms of the sample size of observations, that give rise to the kinetic parameters, getting large. In this context, the above approximation holds if the absolute values of the estimates of the changes in free energy are large compared to their standard errors (see for example the discussions in Hinkley [2] and Marsaglia [3]).

To use this result to derive a distribution for $\widehat{\Phi}$, consider $f(\boldsymbol{\theta}) = f(\Delta\Delta G_{\ddagger}, \Delta\Delta G_U) = \Delta\Delta G_{\ddagger} / \Delta\Delta G_U$. Since

$$\frac{\delta f}{\delta \Delta\Delta G_{\ddagger}} = \frac{1}{\Delta\Delta G_U} \quad \text{and} \quad \frac{\delta f}{\delta \Delta\Delta G_U} = -\frac{\Delta\Delta G_{\ddagger}}{\Delta\Delta G_U^2}, \quad (8)$$

it follows that

$$\widehat{\Phi} = \frac{\widehat{\Delta\Delta G_{\ddagger}}}{\widehat{\Delta\Delta G_U}} \approx N\left(\Phi, \text{var}(\widehat{\Phi})\right), \quad (9)$$

and the variance of the estimate of Φ is given by

$$\begin{aligned} \text{var}(\widehat{\Phi}) &= \frac{1}{(\Delta\Delta G_U)^2} [\Delta\Delta G_U, -\Delta\Delta G_{\ddagger}] \begin{bmatrix} \sigma_{\ddagger}^2 & \rho_{\Delta} \sigma_{\ddagger} \sigma_U \\ \rho_{\Delta} \sigma_{\ddagger} \sigma_U & \sigma_U^2 \end{bmatrix} \begin{bmatrix} \Delta\Delta G_U \\ -\Delta\Delta G_{\ddagger} \end{bmatrix} \frac{1}{(\Delta\Delta G_U)^2} \\ &= \frac{1}{(\Delta\Delta G_U)^4} (\sigma_{\ddagger}^2 \times (\Delta\Delta G_U)^2 - 2\rho_{\Delta} \sigma_{\ddagger} \sigma_U \times \Delta\Delta G_{\ddagger} \Delta\Delta G_U + \sigma_U^2 \times (\Delta\Delta G_{\ddagger})^2). \\ &= \left(\frac{\Delta\Delta G_{\ddagger}}{\Delta\Delta G_U}\right)^2 \times \left[\frac{\sigma_{\ddagger}^2}{(\Delta\Delta G_{\ddagger})^2} - 2\rho_{\Delta} \frac{\sigma_{\ddagger} \sigma_U}{\Delta\Delta G_{\ddagger} \Delta\Delta G_U} + \frac{\sigma_U^2}{(\Delta\Delta G_U)^2} \right] \\ &= \Phi^2 \times \left[\left(\frac{\sigma_{\ddagger}}{\Delta\Delta G_{\ddagger}}\right)^2 - 2\rho_{\Delta} \left(\frac{\sigma_{\ddagger}}{\Delta\Delta G_{\ddagger}}\right) \left(\frac{\sigma_U}{\Delta\Delta G_U}\right) + \left(\frac{\sigma_U}{\Delta\Delta G_U}\right)^2 \right] \end{aligned} \quad (10)$$

2. A Web based Implementation

The above algorithm has been implemented as part of a web server for use by the folding community, accessible through biostat.jhsph.edu/~iruczins/software/phi/. A snapshot of the website is shown in Figure 3. Chevron data can be uploaded as a spreadsheet or as tab delimited text file (there is no limit on the number of mutants that can be analyzed simultaneously). Several user options are available, including the possibility of fitting parallel chevron arms, and the option to measure the folding and unfolding rates at non-zero denaturant concentrations. Once executed, the script creates a webpage with tabulations of the relevant kinetic parameters and their standard errors, as well as estimates of the changes in free energy between mutants, and the estimate for Φ (including the standard error and the confidence interval). These tables can also be downloaded in spreadsheet format.

Figure 3: A snapshot of the web server. A perl/cgi script uploads and parses the kinetic data, and calls an R script for the statistical analysis. The estimates for the kinetic parameters and their standard errors and correlations are obtained using the R function `nls()`.

3. Effects of Ignoring the Covariance

When ρ_{Δ} is ignored (i. e. assumed to be zero), the estimate for the variability in $\hat{\Phi}$ is simply

$$\text{var}_{\rho_{\Delta}=0}(\hat{\Phi}) = \Phi^2 \times \left[\left(\frac{\sigma_{\ddagger}}{\Delta\Delta G_{\ddagger}} \right)^2 + \left(\frac{\sigma_U}{\Delta\Delta G_U} \right)^2 \right]. \quad (11)$$

The absolute difference in those variabilities is

$$\text{var}_{\rho_{\Delta}=0}(\hat{\Phi}) - \text{var}(\hat{\Phi}) = \Phi^2 \times 2\rho_{\Delta} \left(\frac{\sigma_{\ddagger}}{\Delta\Delta G_{\ddagger}} \right) \left(\frac{\sigma_U}{\Delta\Delta G_U} \right) = 2\Phi \times \frac{\rho_{\Delta}\sigma_{\ddagger}\sigma_U}{(\Delta\Delta G_U)^2}. \quad (12)$$

This formula states that the smaller ρ_Δ , σ_\ddagger , σ_u , and Φ (for a given $\Delta\Delta G_u$), the smaller the absolute difference in the estimated variabilities. Also, the larger $\Delta\Delta G_u$ (for a given Φ), the smaller that difference.

To investigate the relative increase in the estimate of the variability of $\hat{\Phi}$ when ρ_Δ is ignored, we simplify the notation by writing

$$\begin{aligned} V_1 &= \Phi^2 (a^2 - 2\rho ab + b^2), \\ V_2 &= \Phi^2 (a^2 + b^2), \end{aligned} \quad (13)$$

where $a = \sigma_\ddagger/\Delta\Delta G_\ddagger$ and $b = \sigma_u/\Delta\Delta G_u$. Using $b = c \times a$, we define

$$\frac{V_1}{V_2} = 1 - \frac{2\rho_\Delta ab}{a^2 + b^2} = 1 - \frac{2\rho_\Delta ca^2}{a^2 + c^2 a^2} = 1 - \frac{2\rho_\Delta c}{1 + c^2} =: f(c) \quad (14)$$

The first derivative of this function f is

$$f'(c) = -\frac{2\rho_\Delta(1-c)}{(1+c^2)^2}, \quad (15)$$

which is equal to zero iff $c = 1$. The function f attains a global minimum at $c = 1$, i.e. when $\sigma_\ddagger/\Delta\Delta G_\ddagger = \sigma_u/\Delta\Delta G_u$. Hence V_1/V_2 is bounded below by $f(1) = 1 - \rho_\Delta$, and converges to 1 (i. e. towards equal variance estimates) as c moves away from this point in either direction. Equivalently, V_2/V_1 is bounded above by $1 - \rho_\Delta/(1 - \rho_\Delta)$. The lengths of the confidence intervals scales with the standard errors of the free energy estimates, and therefore the ratio of confidence intervals is bounded by $\sqrt{V_2}/\sqrt{V_1} = \sqrt{V_2/V_1} = \sqrt{1 - \rho_\Delta/(1 - \rho_\Delta)}$.

References

- [1] Seber, G.A. (1977), *Linear Regression Analysis*, Wiley.
- [2] Hinkley, D.V. (1969), *On the Ratio of Two Correlated Normal Random Variables*, *Biometrika*, 56, 635–39.
- [3] Marsaglia, G. (1965), *Ratios of normal variables and ratios of sums of uniform variables*, *Journal of the American Statistical Association* 60, 193–204.