

Algorithm for enrichment of equilibrium preferences

Jesse D. Bloom^{1*}

¹Division of Basic Sciences and Computational Biology Program,
Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA

*To whom correspondence should be addressed; E-mail: jbloom@fhcrc.org.

Quantifying the preference for amino acids at different sites

Conceptually, our approach is based on the idea that there is an inherent preference for each amino acid at each site in the protein. For now, we assume that these preferences are entirely at the amino-acid level and are indifferent to the specific codon (this assumption is probably not completely accurate, and the study of differential preferences for synonymous codons is an interesting area for future work). We denote the preference of site r for amino-acid a as $\pi_{r,a}$, where these equilibrium preferences are subject to the constraint

$$\sum_a \pi_{r,a} = 1, \quad (1)$$

and where a can range over the $n_{\text{aa}} = 21$ values corresponding to all amino acids plus a stop codon (although in general we expect stop codons to be highly disfavored within the protein). We define the ratio $\pi_{r,a}/\pi_{r,a'}$ to be the expected ratio of amino-acid a to amino-acid a' after viral growth if both are initially introduced into the mutant library at equal frequency. Mutations that enhance viral growth will have larger values of $\pi_{r,a}$, while mutations that hamper viral growth will have lower values of $\pi_{r,a}$. However, note that $\pi_{r,a}/\pi_{r,a'}$ cannot be simply interpreted as the fitness effect of mutating site r from a to a' : because most clones in our mutant libraries have multiple mutations, this ratio summarizes the effect of a mutation in the wildtype gene and a variety of closely related mutants. A mutation can therefore have a ratio greater than one due to its inherent effect on viral growth or its effect on the tolerance of the protein for other mutations (or a combination of both). The analysis here does not separate these two factors, but note that previous experimental work has shown that it is fairly common for one mutation in NP to alter the protein's tolerance to a subsequent mutation.

The most naive approach is to set $\pi_{r,a}$ proportional to the frequency of amino-acid a in the **mutvirus-p1** library divided by the frequency of the mutation in the **mutDNA** library, and then

apply the normalization condition in Equation 1 to get the proportionality constant. However, such an approach is problematic for several reasons. First, it fails to account for other sources of error and mutation (PCR, reverse-transcription, etc) that inflate the observed frequencies of some mutations. Second, the libraries contain finite numbers of counts, and estimating ratios by dividing counts from finite samples is a notoriously statistically biased approach. For example, in the limiting case where a mutation is counted once in the **mutvirus-p1** library and not at all in the **mutDNA** library, the ratio is infinity – yet in practice such low counts give us little confidence that the enough variants have been assayed to estimate the true effect of the mutation.

To circumvent these problems, we use a Bayesian approach that explicitly accounts for the sampling statistics. We begin with prior estimates that the error and mutation rates for each site are equal to the library averages. We specify likelihood functions that give the probability of observing a set of counts given the $\pi_{r,a}$ values for that site and the various error and mutation rates. We then estimate the posterior distribution of the $\pi_{r,a}$ values by MCMC. This approach accounts for sources of error and avoids overfitting $\pi_{r,a}$ when the number counts is low.

We use the counts in the **DNA** library to quantify errors due to PCR and sequencing. We use the counts in the **RNA** library to quantify errors due to reverse-transcription. We assume that transcription of the viral genes from the reverse-genetics plasmids and subsequent replication of these genes by the influenza polymerase introduces a negligible number of new mutations relative to the number already present in the plasmid mutant library. The second of these assumptions is supported by the fact that the total mutation frequency in the **virus-p1** libraries is close to that in the **RNA** libraries. The first of these assumptions is supported by the fact that stop codons are no more frequent in the **RNA** libraries than in the **virus-p1** libraries – deleterious stop codons arising during transcription will be purged during viral growth, while those arising from reverse-transcription and sequencing errors will not.

At each codon site r , there are $n_{\text{codon}} = 64$ codons, which we index by $i = 1, 2, \dots, n_{\text{codon}}$. Let $\text{wt}(r)$ denote the wildtype codon at site r . Let N_r^{DNA} be the total number of sequencing reads at site r in the **DNA** library, and let $n_{r,i}^{\text{DNA}}$ be the number of these reads that report codon i at site r , so that $\sum_i n_{r,i}^{\text{DNA}} = N_r^{\text{DNA}}$. Similarly, let N_r^{mutDNA} , N_r^{RNA} , and N_r^{mutvirus} be the total number of reads at site r and let $n_{r,i}^{\text{mutDNA}}$, $n_{r,i}^{\text{RNA}}$, and $n_{r,i}^{\text{mutvirus}}$ be the total number of these reads that report codon i at site r in the **mutDNA**, **RNA**, and **mutvirus-p1**, respectively.

We first consider the rate at which site r is erroneously read to be some incorrect identity due to PCR or sequencing errors. Such errors are the only source of non-wildtype reads in the sequencing of the **DNA** library. For all $i \neq \text{wt}(r)$, we define $\epsilon_{r,i}$ as the rate at which site r is erroneously read as codon i in the **DNA** library. We define $\epsilon_{r,\text{wt}(r)} = 1 - \sum_{i \neq \text{wt}(r)} \epsilon_{r,i}$ to be

the rate at which site r is correctly read as its wildtype identity of $\text{wt}(r)$ in the **DNA** library. We therefore have $\epsilon_{r,i} = \mathbb{E} [n_{r,i}^{\text{DNA}} / N_r^{\text{DNA}}]$ where \mathbb{E} denotes the expectation value. If we define $\vec{\epsilon}_r = (\epsilon_{r,1}, \dots, \epsilon_{r,n_{\text{codon}}})$ and $\vec{n}_r^{\text{DNA}} = (n_{r,1}^{\text{DNA}}, \dots, n_{r,n_{\text{codon}}}^{\text{DNA}})$ as vectors of the $\epsilon_{r,i}$ and $n_{r,i}^{\text{DNA}}$ values, then the likelihood of observing \vec{n}_r^{DNA} given $\vec{\epsilon}_r$ and N_r^{DNA} is

$$\Pr \left(\vec{n}_r^{\text{DNA}} \mid N_r^{\text{DNA}}, \vec{\epsilon}_r \right) = \text{Mult} \left(\vec{n}_r^{\text{DNA}}; N_r^{\text{DNA}}, \vec{\epsilon}_r \right) \quad (2)$$

where Mult denotes the multinomial distribution.

We next consider the rate at which site r is erroneously copied during reverse transcription. These reverse-transcription errors combine with the PCR / sequencing errors defined by $\vec{\epsilon}_r$ to create non-wildtype reads in the **RNA** library. For all $i \neq \text{wt}(r)$, we define $\rho_{r,i}$ as the rate at which site r is miscopied to codon i during reverse transcription. We define $\rho_{r,\text{wt}(r)} = 1 - \sum_{i \neq \text{wt}(r)} \rho_{r,i}$ as the rate at which site r is correctly reverse transcribed. If we ignore as negligibly rare the possibility that a site is subject to both a reverse-transcription and sequencing / PCR error within the same clone (a reasonable assumption as both $\epsilon_{r,i}$ and $\rho_{r,i}$ are very small for $i \neq \text{wt}(r)$), then $\epsilon_{r,i} + \rho_{r,i} - \delta_{i,\text{wt}(r)} = \mathbb{E} [n_{r,i}^{\text{RNA}} / N_r^{\text{RNA}}]$ where $\delta_{i,\text{wt}(r)}$ is the Kronecker delta (equal to one if $i = \text{wt}(r)$ and zero otherwise). The likelihood of observing \vec{n}_r^{RNA} given $\vec{\rho}_r, \vec{\epsilon}_r$, and N_r^{RNA} is

$$\Pr \left(\vec{n}_r^{\text{RNA}} \mid N_r^{\text{RNA}}, \vec{\rho}_r, \vec{\epsilon}_r \right) = \text{Mult} \left(\vec{n}_r^{\text{RNA}}; N_r^{\text{RNA}}, \vec{\epsilon}_r + \vec{\rho}_r - \vec{\delta}_r \right). \quad (3)$$

where $\vec{\delta}_r = (\delta_{1,\text{wt}(r)}, \dots, \delta_{n_{\text{codon}},\text{wt}(r)})$ is a vector that is all zeros except for the element corresponding to $\text{wt}(r)$.

We next consider the rate at which site r is mutated to some other codon in the plasmid mutant library. These mutations combine with the PCR / sequencing errors defined by $\vec{\epsilon}_r$ to create non-wildtype reads in the **mutDNA** library. For all $i \neq \text{wt}(r)$, we define $\mu_{r,i}$ as the rate at which site r is mutated to codon i in the mutant library. We define $\mu_{r,\text{wt}(r)} = 1 - \sum_{i \neq \text{wt}(r)} \mu_{r,i}$ as the rate at which site r is not mutated. If we ignore as negligibly rare the possibility that a site is subject to both a mutation and a sequencing / PCR error within the same clone, then $\mu_{r,i} + \epsilon_{r,i} - \delta_{i,\text{wt}(r)} = \mathbb{E} [n_{r,i}^{\text{mutDNA}} / N_r^{\text{mutDNA}}]$. The likelihood of observing $\vec{n}_r^{\text{mutDNA}}$ given $\vec{\mu}_r, \vec{\epsilon}_r$, and N_r^{mutDNA} is

$$\Pr \left(\vec{n}_r^{\text{mutDNA}} \mid N_r^{\text{mutDNA}}, \vec{\mu}_r, \vec{\epsilon}_r \right) = \text{Mult} \left(\vec{n}_r^{\text{mutDNA}}; N_r^{\text{mutDNA}}, \vec{\mu}_r + \vec{\epsilon}_r - \vec{\delta}_r \right). \quad (4)$$

Finally, we consider the effect of the preferences of each site r for different amino acids, as denoted by the $\pi_{r,a}$ values. Selection due to these preferences is manifested in the **mutvirus** library. This selection acts on the mutations in the mutant library ($\mu_{r,i}$), although the actual counts in the **mutvirus** library are also affected by the sequencing / PCR errors ($\epsilon_{r,i}$) and reverse-transcription errors ($\rho_{r,i}$). We again ignore as negligibly rare the possibility that a site is subject to more than one of these sources of mutation and error within a single clone. Let $\mathcal{A}(i)$ denote the amino acid encoded by codon i . Let $\vec{\pi}_r$ be the vector of $\pi_{r,a}$ values. Define the vector-valued function $\vec{\mathcal{C}}$ as

$$\vec{\mathcal{C}}(\vec{\pi}_r) = (\pi_{r,\mathcal{A}(1)}, \dots, \pi_{r,\mathcal{A}(n_{\text{codon}})}), \quad (5)$$

so that this function returns a n_{codon} -element vector constructed from $\vec{\pi}_r$. Because the selection in the **mutvirus** library due to the preferences $\pi_{r,\mathcal{A}(i)}$ occurs after the mutagenesis $\mu_{r,i}$ but before the reverse-transcription errors $\rho_{r,i}$ and the sequencing / PCR errors $\epsilon_{r,i}$, we have

$\mathbb{E} \left[n_{r,i}^{\text{mutvirus}} / N_r^{\text{mutvirus}} \right] = \epsilon_{r,i} + \rho_{r,i} + \gamma_r \times \pi_{r,\mathcal{A}(i)} \times \mu_{r,i} - 2 \times \delta_{i,\text{wt}(r)}$ where $\gamma_r = \left(\sum_i \pi_{r,\mathcal{A}(i)} \mu_{r,i} \right)^{-1}$ (where $\vec{\mathcal{C}}(\vec{\pi}_r) \cdot \vec{\mu}_r$)⁻¹ (where \cdot denotes the dot product) is a normalization factor that accounts for the fact that changes in the frequency of one variant due to selection will influence the observed frequency of other variants. The likelihood of observing $\overrightarrow{n_r^{\text{mutvirus}}}$ given $\vec{\mu}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\pi}_r$, and N_r^{mutvirus} is therefore

$$\Pr \left(\overrightarrow{n_r^{\text{mutvirus}}} \mid \vec{\mu}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\pi}_r, N_r^{\text{mutvirus}} \right) = \text{Mult} \left(\overrightarrow{n_r^{\text{mutvirus}}}; N_r^{\text{mutvirus}}, \vec{\epsilon}_r + \vec{\rho}_r + \frac{\vec{\mathcal{C}}(\vec{\pi}_r) \circ \vec{\mu}_r}{\vec{\mathcal{C}}(\vec{\pi}_r) \cdot \vec{\mu}_r} - 2\vec{\delta}_r \right). \quad (6)$$

where \circ is the Hademard (entry-wise) product.

We specify priors over $\vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r$, and $\vec{\pi}_r$ in the form of Dirichlet distributions (denoted here by Dir). For the priors over the mutation rates $\vec{\mu}_r$, we choose Dirichlet-distribution parameters such that the mean of the prior expectation for the mutation rate at each site r and codon i is proportional to the average value for all sites. The average fraction of mutated codons in the **mutDNA** library minus the background from the **DNA** library is 6.1×10^{-3} , so the average mutation rate is $\bar{\mu} = 6.1 \times 10^{-3} / 63 = 9.7 \times 10^{-5}$. So we use a prior of

$$\Pr(\vec{\mu}_r) = \text{Dir}(\vec{\mu}_r; n_{\text{codon}} \cdot \sigma_\mu \cdot \overrightarrow{\alpha_{\mu,r}}) \quad (7)$$

where $\overrightarrow{\alpha_{\mu,r}}$ is the n_{codon} -element vector with elements $\alpha_{\mu,r,i} = \bar{\mu} + \delta_{i,\text{wt}(r)} (1 - n_{\text{codon}} \bar{\mu})$ and σ_μ is the scalar concentration parameter. For a symmetric Dirichlet distribution (which $\Pr(\vec{\mu}_r)$ will in general *not* be), choosing $\sigma_\mu = 1$ makes the distribution completely uniform.

For the priors over $\epsilon_{r,i}$ and $\rho_{r,i}$, the Dirichlet-distribution parameters again represent the average value for all sites, but now also depend on the number of nucleotide changes in the codon mutation since sequencing / PCR and reverse-transcription errors are far more likely to lead to single-nucleotide codon changes than multiple-nucleotide codon changes. Let $\mathcal{M}(\text{wt}(r), i)$ be the number of nucleotide changes in the mutation from codon $\text{wt}(r)$ to codon i . For example, $\mathcal{M}(\text{GCA}, \text{ACA}) = 1$ and $\mathcal{M}(\text{GCA}, \text{ATA}) = 2$. The average error rate (estimated from the **DNA** library) is $\bar{\epsilon}_1 = 5.8 \times 10^{-4} / 9 = 6.4 \times 10^{-5}$ for single-nucleotide codon mutations, $\bar{\epsilon}_2 = 8.7 \times 10^{-6} / 27 = 3.2 \times 10^{-7}$ for two-nucleotide codon mutations, and $\bar{\epsilon}_3 = 4.0 \times 10^{-6} / 27 = 1.5 \times 10^{-7}$ for three-nucleotide codon mutations. So we use a prior of

$$\Pr(\vec{\epsilon}_r) = \text{Dir}(\vec{\epsilon}_r; n_{\text{codon}} \cdot \sigma_\epsilon \cdot \overrightarrow{\alpha_{\epsilon,r}}) \quad (8)$$

where $\overrightarrow{\alpha_{\epsilon,r}}$ is the n_{codon} -element vector with elements $\alpha_{\epsilon,r,i} = \overline{\epsilon_{\mathcal{M}(\text{wt}(r), i)}}$ where we define $\bar{\epsilon}_0 = 1 - 9 \times \bar{\epsilon}_1 - 27 \times \bar{\epsilon}_2 - 27 \times \bar{\epsilon}_3$ (for any codon, there are 9 one-nucleotide mutations, 27 two-nucleotide mutations, and 27 three-nucleotide mutations), and where σ_ϵ is the scalar concentration parameter.

Similarly, the average reverse-transcription error rates (estimated from the **RNA** library minus the **DNA** library) are $\bar{\rho}_1 = 1.9 \times 10^{-4} / 9 = 2.1 \times 10^{-5}$, $\bar{\rho}_2 = 1.5 \times 10^{-5} / 27 = 5.6 \times 10^{-7}$,

and $\bar{\rho}_3 = 2.8 \times 10^{-6}/27 = 1.0 \times 10^{-7}$ for one-, two-, and three-nucleotide codon mutations, respectively. So we use a prior of

$$\Pr(\vec{\rho}_r) = \text{Dir}(\vec{\rho}_r; n_{\text{codon}} \cdot \sigma_\rho \cdot \vec{\alpha}_{\rho,r}) \quad (9)$$

where $\vec{\alpha}_{\rho,r}$ is the n_{codon} -element vector with elements $\alpha_{\rho,r,i} = \overline{\rho_{\mathcal{M}(\text{wt}(r),i)}}$ where we define $\bar{\rho}_0 = 1 - 9 \times \bar{\rho}_1 - 27 \times \bar{\rho}_2 - 27 \times \bar{\rho}_3$, and where σ_ρ is the scalar concentration parameter.

We specify a symmetric Dirichlet-distribution prior over $\vec{\pi}_r$ (note that any other prior, such as one that favored the wildtype identity, would implicitly favor certain identities based empirically on the wildtype protein sequence, and so would not be in the spirit of the parameter-free derivation of the $\pi_{r,a}$ values employed here). Specifically, we use a prior of

$$\Pr(\vec{\pi}_r) = \text{Dir}(\vec{\pi}_r; \sigma_\pi \cdot \vec{1}) \quad (10)$$

where $\vec{1}$ is the n_{aa} -element vector that is all ones, and σ_π is the scalar concentration parameter. Setting $\sigma_\pi = 1$ gives a uniform prior over all $\vec{\pi}_r$ values, while smaller values favor peaked distribution and larger values favor equal values for all $\pi_{r,a}$ elements.

We can now write expressions for the likelihoods and posterior probabilities. Let $\mathcal{N}_r = \{n_r^{\text{DNA}}, n_r^{\text{mutDNA}}, n_r^{\text{RNA}}, n_r^{\text{mutvirus}}, N_r^{\text{DNA}}, N_r^{\text{mutDNA}}, N_r^{\text{RNA}}, N_r^{\text{mutvirus}}\}$ denote the full set of counts for site r . The likelihood of \mathcal{N}_r given values for the equilibrium preferences and mutation / error rates is

$$\begin{aligned} \Pr(\mathcal{N}_r | \vec{\pi}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r) &= \Pr(n_r^{\text{DNA}} | N_r^{\text{DNA}}, \vec{\epsilon}_r) \times \Pr(n_r^{\text{RNA}} | N_r^{\text{RNA}}, \vec{\epsilon}_r, \vec{\rho}_r) \times \\ &\Pr(n_r^{\text{mutDNA}} | N_r^{\text{mutDNA}}, \vec{\epsilon}_r, \vec{\mu}_r) \times \\ &\Pr(n_r^{\text{mutvirus}} | N_r^{\text{mutvirus}}, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r, \vec{\pi}_r) \end{aligned} \quad (11)$$

where the likelihoods that compose Equation 11 are defined by Equations 2, 3, 4, and 6. The posterior probability of a specific value for the equilibrium preferences and mutation / error rates is

$$\Pr(\vec{\pi}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r | \mathcal{N}_r) = C_r \times \Pr(\mathcal{N}_r | \vec{\pi}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r) \times \quad (12)$$

$$\Pr(\vec{\epsilon}_r) \times \Pr(\vec{\rho}_r) \times \Pr(\vec{\mu}_r) \times \Pr(\vec{\pi}_r) \quad (13)$$

where C_r is a normalization constant that does not need to be explicitly calculated in the MCMC approach used here. The posterior over the equilibrium preferences $\vec{\pi}_r$ can be calculated by integrating over Equation 12 to give

$$\Pr(\vec{\pi}_r | \mathcal{N}_r) = \int \int \int \Pr(\vec{\pi}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r | \mathcal{N}_r) d\vec{\epsilon}_r d\vec{\rho}_r d\vec{\mu}_r, \quad (14)$$

where the integration is performed by MCMC.

Equation 14 infers $\vec{\pi}_r$ for a single replicate of the experiment. In practice, we have performed $\mathcal{R} = 4$ replicates (WT-1, WT-2, N334H-1, N334H-2). Let \mathcal{N}_r^k denote the set of experimentally observed counts for replicate k . The full set of data for all \mathcal{R} replicates is then $\{\mathcal{N}_r^k \mid 1 \leq k \leq \mathcal{R}\}$. We use the same equilibrium preferences $\vec{\pi}_r$ for all replicates, since these preferences are a property of the protein rather than the experiment. However, the error and mutation rates $\vec{\epsilon}_r$, $\vec{\mu}_r$, and $\vec{\rho}_r$ are specific to each replicate (denote the replicate-specific vectors by $\vec{\epsilon}_r^k$, $\vec{\mu}_r^k$, and $\vec{\rho}_r^k$), with the replicate-specific prior vectors $\vec{\alpha}_{\mu,r}^k$, $\vec{\alpha}_{\epsilon,r}^k$, and $\vec{\alpha}_{\rho,r}^k$ determined from the library averages for that replicate. The reason for making these three vectors and their priors replicate specific is that each replicate could have different error and mutation rates. The posterior probability of $\vec{\pi}_r$ and the mutation / error rates for all experimental replicates is

$$\Pr(\vec{\pi}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r \mid \{\mathcal{N}_r^k \mid 1 \leq k \leq \mathcal{R}\}) = C'_r \times \prod_{k=1}^{\mathcal{R}} \left[\Pr(\mathcal{N}_r^k \mid \vec{\pi}_r, \vec{\epsilon}_r^k, \vec{\rho}_r^k, \vec{\mu}_r^k) \times \Pr(\vec{\epsilon}_r^k) \Pr(\vec{\rho}_r^k) \Pr(\vec{\mu}_r^k) \right] \Pr(\vec{\pi}_r), \quad (15)$$

where each of the individual likelihoods for the \mathcal{N}_r^k counts are calculated using Equation 11 and C'_r is a normalization constant that does not need to be calculated when using MCMC. The posterior over the equilibrium preferences $\vec{\pi}_r$ can be calculated by using MCMC to integrate over Equation 15 to give

$$\Pr(\vec{\pi}_r \mid \{\mathcal{N}_r^k \mid 1 \leq k \leq \mathcal{R}\}) = \int \cdots \int \Pr(\vec{\pi}_r, \vec{\epsilon}_r, \vec{\rho}_r, \vec{\mu}_r \mid \{\mathcal{N}_r^k \mid 1 \leq k \leq \mathcal{R}\}) \prod_{k=1}^{\mathcal{R}} d\vec{\epsilon}_r^k d\vec{\rho}_r^k d\vec{\mu}_r^k. \quad (16)$$

We summarize the posterior calculated from Equation 16 by its mean,

$$\langle \vec{\pi}_r \rangle = \int \vec{\pi}_r \times \Pr(\vec{\pi}_r \mid \{\mathcal{N}_r^k \mid 1 \leq k \leq \mathcal{R}\}) d\vec{\pi}_r. \quad (17)$$