

TnSeq Data Analysis

Sergey Kryazhimskiy

November 29, 2017

Pilot 2017-06

Identifying selected mutations (simple model, no reads)

We will model the selection coefficient X_r for replicate $r = 1, 2, \dots, n$ of a particular transformant as

$$X_r = \mu + \delta + \varepsilon_r,$$

where μ is the true effect of the transposon-insertion mutation, δ is a defect potentially introduced by transformation, and ε_r is the measurement error. We make the following assumptions about these variables. μ is a parameter (not a random variable). ε_r are all i.i.d r.v.'s, normally distributed with mean 0 and variance σ_{err}^2 . δ is an r.v. with a mixture distribution, such that

$$\delta = \begin{cases} 0, & \text{with probability } 1 - p_{\text{tr}}, \\ -\Delta, & \text{with probability } p_{\text{tr}}, \end{cases}$$

where Δ are normally distributed with parameters μ_{tr} and σ_{tr}^2 . In other words, with probability p_{tr} , transformation introduces an artifact, which is a normally distributed (typically deleterious) effect; and with probability $1 - p_{\text{tr}}$, no such artifact is introduced. The conditional probability of observing the selection coefficient X_r in the vicinity of value x_r in replicate experiment r , given that the current transformant's true selection coefficient Y in the given environment, is

$$\mathbb{P}_{X_r|Y}(dx_r) \equiv p(x_r|Y) dx_r = N_1(x_r; Y, \sigma_{\text{err}}^2) dx_r,$$

where $N_d(\mathbf{x}; \mathbf{m}, \mathbf{\Sigma})$ is a d -dimensional multivariate gaussian probability density function with respect to variable \mathbf{x} , with mean \mathbf{m} and variance-covariance matrix $\mathbf{\Sigma}$. The joint conditional probability density that the selection coefficients $\mathbf{X} = (X_1, \dots, X_n)$ in replicate measurements of a particular transformant in a particular environment are in the vicinity of point $\mathbf{x} = (x_1, \dots, x_n)$, given the transformant's true fitness Y in that environment, is

$$p(\mathbf{x}|Y; \sigma_{\text{err}}) = \prod_{r=1}^n N_1(x_r; Y, \sigma_{\text{err}}) = \frac{\exp\left\{-\frac{(n-1)\text{Var } \mathbf{x}}{2\sigma_{\text{err}}^2}\right\}}{\sqrt{n(2\pi\sigma_{\text{err}}^2)^{n-1}}} N_1\left(Y; \bar{x}, \frac{\sigma_{\text{err}}}{\sqrt{n}}\right), \quad (1)$$

where $\bar{x} = \frac{1}{n} \sum_{r=1}^n x_r$ and $\text{Var } \mathbf{x} = \frac{1}{n-1} \sum_{r=1}^n (x_r - \bar{x})^2$.

Next, we will assume that a transformation artifact occurs with probability p_{tr} . If it occurs, its effect $\mathbf{Y} = (Y_1, \dots, Y_K)$ in environments e_1, e_2, \dots, e_K is drawn from a multivariate normal distribution with mean vector $\boldsymbol{\mu}_{\text{tr}} > 0$ and variance-covariance matrix $\boldsymbol{\Sigma}_{\text{tr}}$. Suppose that the true selection coefficients of the current transformant in environments e_1, e_2, \dots, e_K are given by $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$. Then,

$$\mathbb{P}_{\mathbf{Y}}(d\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) = (1 - p_{\text{tr}}) \delta_{\boldsymbol{\mu}}(d\mathbf{y}) + p_{\text{tr}} N_K(\mathbf{y}; \boldsymbol{\mu} - \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) d\mathbf{y},$$

where $\delta_{\boldsymbol{\mu}}(d\mathbf{y})$ is a point measure at $\boldsymbol{\mu}$.

Now consider all replicate measurements of the same transformant in all environments. We have $\mathbf{X}_1 = (X_{11}, \dots, X_{1n_1})$ measurements of the selection coefficient of the transformant in environment 1, $\mathbf{X}_2 = (X_{21}, \dots, X_{2n_2})$ measurements in environment 2, etc. up to environment K . The joint probability density of such observation is given by

$$\begin{aligned} & p(\mathbf{x}_1, \dots, \mathbf{x}_K; \boldsymbol{\mu}, \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}, \boldsymbol{\sigma}_{\text{err}}) \\ &= \int_{\mathbb{R}^K} p(\mathbf{x}_1, \dots, \mathbf{x}_K | \mathbf{Y}; \boldsymbol{\sigma}_{\text{err}}) \mathbb{P}_{\mathbf{Y}}(d\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) \\ &= \int_{\mathbb{R}^K} \prod_{k=1}^K p(\mathbf{x}_k | Y_k; \sigma_{\text{err},k}) [(1 - p_{\text{tr}}) \delta_{\boldsymbol{\mu}}(d\mathbf{y}) + p_{\text{tr}} N_K(\mathbf{y}; \boldsymbol{\mu} - \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) d\mathbf{y}] \\ &= \int_{\mathbb{R}^K} \prod_{k=1}^K \frac{\exp\left\{-\frac{(n_{\text{rep},k} - 1) \text{Var } \mathbf{x}_k}{2\sigma_{\text{err},k}^2}\right\}}{\sqrt{n_{\text{rep},k} (2\pi\sigma_{\text{err},k}^2)^{n_{\text{rep},k} - 1}}} N_1\left(y_k; \bar{x}_k, \frac{\sigma_{\text{err}}}{\sqrt{n_{\text{rep},k}}}\right) \\ &\times [(1 - p_{\text{tr}}) \delta_{\boldsymbol{\mu}}(d\mathbf{y}) + p_{\text{tr}} N_K(\mathbf{y}; \boldsymbol{\mu} - \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) d\mathbf{y}] \\ &= \frac{\exp\left\{-\sum_{k=1}^K (n_{\text{rep},k} - 1) \frac{\text{Var } \mathbf{x}_k}{2\sigma_{\text{err},k}^2}\right\}}{\sqrt{\prod_{k=1}^K n_{\text{rep},k} (2\pi\sigma_{\text{err},k}^2)^{n_{\text{rep},k} - 1}}} \int_{\mathbb{R}^K} N_K(\mathbf{y}; \mathbf{M}, \boldsymbol{\Sigma}_{\text{err}}) \\ &\times [(1 - p_{\text{tr}}) \delta_{\boldsymbol{\mu}}(d\mathbf{y}) + p_{\text{tr}} N_K(\mathbf{y}; \boldsymbol{\mu} - \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) d\mathbf{y}] \\ &= \frac{\exp\left\{-\sum_{k=1}^K (n_{\text{rep},k} - 1) \frac{\text{Var } \mathbf{x}_k}{2\sigma_{\text{err},k}^2}\right\}}{\sqrt{\prod_{k=1}^K n_{\text{rep},k} (2\pi\sigma_{\text{err},k}^2)^{n_{\text{rep},k} - 1}}} \left[(1 - p_{\text{tr}}) N_K(\boldsymbol{\mu}; \mathbf{M}, \boldsymbol{\Sigma}_{\text{err}}) \right. \\ &\left. + p_{\text{tr}} \int_{\mathbb{R}^K} N_K(\mathbf{y}; \mathbf{M}, \boldsymbol{\Sigma}_{\text{err}}) N_K(\mathbf{y}; \boldsymbol{\mu} - \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}) d\mathbf{y} \right] \end{aligned} \quad (2)$$

where to obtain the last equality we used relationship (1). In equation (2), $\mathbf{M} = (\bar{x}_1, \dots, \bar{x}_K)^T$ and $\boldsymbol{\Sigma}_{\text{err}}$ is a diagonal matrix with $(\sigma_{\text{err}}/\sqrt{n_{\text{rep},1}}, \dots, \sigma_{\text{err}}/\sqrt{n_{\text{rep},K}})$ on the diagonal, with $\bar{x}_k = \frac{1}{n_{\text{rep},k}} \sum_{r=1}^{n_{\text{rep},k}} x_{kr}$ and $\text{Var } \mathbf{x}_k = \frac{1}{n_{\text{rep},k} - 1} \sum_{r=1}^{n_{\text{rep},k}} (x_{kr} - \bar{x}_k)^2$. Using properties of the multivariate

normal distribution, the integral in equation (2) can be taken analytically, which gives us

$$\begin{aligned}
& p(\mathbf{x}_1, \dots, \mathbf{x}_K; \boldsymbol{\mu}, \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{tr}}, \boldsymbol{\sigma}_{\text{err}}) \\
&= \frac{\exp\left\{-\sum_{k=1}^K (n_{\text{rep},k} - 1) \frac{\text{Var } \mathbf{x}_k}{2\sigma_{\text{err},k}^2}\right\}}{\sqrt{\prod_{k=1}^K n_{\text{rep},k} (2\pi\sigma_{\text{err},k}^2)^{n_{\text{rep},k}-1}}} \\
&\times \left[(1 - p_{\text{tr}}) N_K(\mathbf{M}; \boldsymbol{\mu}, \boldsymbol{\Sigma}_{\text{err}}) + p_{\text{tr}} N_K(\mathbf{M}; \boldsymbol{\mu} - \boldsymbol{\mu}_{\text{tr}}, \boldsymbol{\Sigma}_{\text{err}} + \boldsymbol{\Sigma}_{\text{tr}}) \right] \quad (3)
\end{aligned}$$

Identifying fraction and mean effect of selected mutations

To identify the fraction of beneficial mutations p^b , fraction of deleterious mutations p^n , and the fraction of deleterious mutations p^d , as well as the mean effect of deleterious and beneficial mutations, we will employ the following model. We will assume that, for each mutation m in strain s and environment e , the effect comes from a mixture distribution:

$$\mu_{sem} = \begin{cases} 0, & \text{with probability } p_{n,se}, \\ -X_{d,sem}, & \text{with probability } p_{d,se}, \\ X_{b,sem}, & \text{with probability } p_{b,se}, \end{cases}$$

where $X_{b,sem}$ and $X_{d,sem}$ are exponentially distributed random variables with parameters $\theta_{d,se}$ and $\theta_{b,se}$, respectively. Thus, the probability of observing the selection coefficient X_{selbr} is given by

$$\begin{aligned}
& \mathbb{P}(X_{selbr}; \mathbf{p}_{se}, \boldsymbol{\theta}_{se}, p_{\text{tr},s}, \mu_{\text{tr},se}, \sigma_{\text{tr},se}^2, \sigma_{\text{err},e}^2) = p_{n,se} \mathbb{P}(X_{selbr}; 0, p_{\text{tr},s}, \mu_{\text{tr},se}, \sigma_{\text{tr},se}^2, \sigma_{\text{err},e}^2) \\
&+ p_{se}^d \int_0^\infty \mathbb{P}(X_{selbr}; -x, p_{\text{tr},s}, \mu_{\text{tr},se}, \sigma_{\text{tr},se}^2, \sigma_{\text{err},e}^2) P_{\text{exp}}(x; \theta_{d,se}) dx \\
&+ p_{se}^b \int_0^\infty \mathbb{P}(X_{selbr}; x, p_{\text{tr},s}, \mu_{\text{tr},se}, \sigma_{\text{tr},se}^2, \sigma_{\text{err},e}^2) P_{\text{exp}}(x; \theta_{b,se}) dx \quad (4)
\end{aligned}$$

Note on estimating the error distribution

Suppose we have measurements grouped into n groups $X_{11}, X_{12}, \dots, X_{1m_1}, X_{21}, \dots, X_{nm_n}$. Within group i , all measurements are i.i.d. with mean μ_i . Let us assume that the distributions for each group are the same, except with a shifted mean. Thus, the variance is the same in all groups, σ^2 . Think of each group as replicate measurements of the same barcode, and different groups are measurements of the same mutation but different barcodes. The variance comes only from measurement noise, which is the same for all barcodes.

We would like to obtain a non-parametric estimate of the underlying error distribution. Define \bar{X}_i as the group mean, as usual. Let $Y_{ij} = X_{ij} - \bar{X}_i$. We want to pool all Y_{ij} together to estimate the distribution.

First, we note that $\mathbb{E}Y_{ij} = 0$. Next, we have

$$\begin{aligned}
 \text{Var } Y_{ij} &= \mathbb{E}Y_{ij}^2 - \underbrace{(\mathbb{E}Y_{ij})^2}_{=0} = \mathbb{E}X_{ij}^2 - 2\mathbb{E}(X_{ij}\bar{X}_i) + \mathbb{E}\bar{X}_i^2 \\
 &= \begin{bmatrix} \mathbb{E}X_{ij}^2 = \sigma^2 + \mu_i^2 \\ \mathbb{E}\bar{X}_i^2 = \frac{\sigma^2}{m_i} + \mu_i^2 \\ \mathbb{E}(X_i\bar{X}) = \frac{\sigma^2}{m_i} + \mu_i^2 \end{bmatrix} \\
 &= \sigma^2 \frac{m_i - 1}{m_i}.
 \end{aligned}$$

Thus, variables $\tilde{Y}_{ij} = \sqrt{\frac{m_i}{m_i - 1}}Y_{ij}$ will have the same variance, and hence will be i.i.d.