# 1. Background parameterization

Link two parameterization of gamma distribution

$$gamma(x) = g_1(c,s) = \frac{1}{\Gamma(c)s^c} x^{c-1} e^{-\frac{x}{s}} = \frac{\left(\frac{x}{s}\right)^{c-1} e^{-\frac{x}{s}}}{s\Gamma(c)}$$

$$= g_2(\mu,\sigma) = \frac{1}{(\sigma^2\mu)^{\frac{1}{\sigma^2}}} \frac{y^{\left(\frac{1}{\sigma^2}-1\right)} \exp\left(-\frac{x}{\sigma^2\mu}\right)}{\Gamma\left(\frac{1}{\sigma^2}\right)} \tag{0.1}$$

to link the two function, we have

$$\begin{cases} c = \frac{1}{\sigma^2} \\ s = \sigma^2\mu \end{cases} \Leftrightarrow \begin{cases} \sigma^2 = \frac{1}{c} \\ \mu = sc \end{cases} \tag{0.2}$$

Then we have in general the function:

$$g_1(c,s) = g_2(\mu,\sigma) = g_1\left(\frac{1}{\sigma^2},\sigma^2\mu\right) = g_2\left(sc,\sqrt{\frac{1}{c}}\right) \tag{0.3}$$

Some properties of gamma distribution: convolutions:

$$f_{X_1}(x) = g_1(c_1,s), f_{X_2}(x) = g_1(c_2,s) \Leftrightarrow f_{X_1+X_2}(x) = g_1(c_1+c_2,s) \tag{0.4}$$

$$f_{X_k}(x) = g_1(c_k,s) \Leftrightarrow f_{\sum X_k}(x) = g_1\left(\sum c_k,s\right) \tag{0.5}$$

Equations (0.4) and (0.5) are well known but not directly interpretable. Let set them into mean and variance format as

$$f_{X_k}(x) = g_1(c_k,s) = g_2\left(sc_k,\sqrt{1/c_k}\right) \Leftrightarrow f_{\sum X_k}(x) = g_1\left(\sum c_k,s\right) = g_2\left(s\sum c_k,\sqrt{\frac{1}{\sum c_k}}\right) \tag{0.6}$$

# 2. using Poisson-gamma Mixture to sample the mutation rate heterogeneity among samples

Suppose $y_i^{d,s}$ and $\lambda_i^{d,s}$ denote the mutation count and mutation rate for bin $i$, disease $d$ and sample $s$. $S_d$ represents the number of sample in disease $d$. We assume that the conditional distribution of $y_i^{d,s}$ follows a Poisson distribution with PMF.

$$P\left\{y_i^{d,s}\big|\lambda_i^{d,s}\right\} = \frac{\left(\lambda_i^{d,s}\right)^{y_i^{d,s}} \exp\left\{-\left(\lambda_i^{d,s}\right)\right\}}{\left(y_i^{d,s}\right)!} \tag{1}$$

When these samples are independent, we pool the samples from the same disease by

$$y_i^d = \sum_{s=1}^{S_d} y_i^{d,s} \tag{2}$$

When $\lambda_i^{d,s}$ is fixed (nonrandom but still can be different, and no conditional needed), the PMF of $y_i^d$ can be written into

$$P\left\{y_i^d\right\} = \frac{\left(\sum_{s=1}^{S_d} \lambda_i^{d,s}\right)^{y_i^d} \exp\left\{-\left(\sum_{s=1}^{S_d} \lambda_i^{d,s}\right)\right\}}{\left(y_i^d\right)!} \tag{3}.$$

Specifically, in a constant rate assumption, $\lambda_i^{d,s} \triangleq \lambda$, the equation (3) can be written as

$$P\left\{y_i^d\right\} = \frac{\left(S_d \lambda\right)^{y_i^d} \exp\left(-S_d \lambda\right)}{\left(y_i^d\right)!} \tag{4}.$$

Now in our model, we assume that When $\lambda_i^{d,s}$ i.i.d gamma random variables, and its distribution is

$$p\left(\lambda_i^{d,s} = x\right) = \left(\frac{x}{s}\right)^{c-1} \frac{\exp\left(\dfrac{x}{s}\right)}{\left\{s\Gamma(c)\right\}} = g_1(s,c) \tag{5}.$$

Then we have the distribution of pooled mutation rate as

$$p\left(\lambda_i^d = x\right) = p\left(\sum_{s=1}^{S_d} \lambda_i^{d,s} = x\right) \sim \left(\left(\frac{x}{s}\right)^{nc-1} \frac{\exp\left(\dfrac{x}{s}\right)}{\left\{s\Gamma(nc)\right\}}\right) = g_1(nc,s) \tag{6}.$$

We may rewrite (3) with the $\lambda_i^d$ (random variable) as

$$P\left\{y_i^d \middle| \lambda_i^d\right\} = \frac{\left(\lambda_i^d\right)^{y_i^d} \exp\left(\lambda_i^d\right)}{\left(\lambda_i^d\right)!} \tag{7}.$$

Putting (6) & (7) together we have

$$P\left(y_i^d = y\right) = \left(\frac{1}{1+s}\right)^{nc} \frac{\Gamma(nc+y)}{\Gamma(nc)y!}\left(\frac{s}{1+s}\right)^y \tag{8}.$$

Then the mean and variance can be expressed by $E(y) = nsc$ and $\operatorname{var}(y) = nsc(1+s) = (1+s)E(y)$.

Let

$$\left\{\begin{array}{l} 1/\sigma' = nc \\ s = \sigma'\mu' \end{array}\right. \Leftrightarrow \left\{\begin{array}{l} c = 1/(n\sigma') \\ s = \sigma'\mu' \end{array}\right. \Leftrightarrow \left\{\begin{array}{l} \sigma' = 1/(nc) \\ \mu' = nsc \end{array}\right. \tag{9}$$

Put (9) into (8), we can re-parameterize our NBI distribution using $\acute{\mu}$ and $\acute{\sigma}$ notations as

$$p\left(y_i^d = y \middle| \mu',\sigma'\right) = \frac{\Gamma\left(y+1/\sigma'\right)}{\Gamma\left(1/\sigma'\right)y!}\left(\frac{\sigma'\mu'}{1+\sigma'\mu'}\right)^y\left(\frac{1}{1+\sigma'\mu'}\right)^{(1/\sigma')} \tag{10}.$$

It can be regarded as a Poisson-gamma mixture distribution with

$$P(Y|\mu\gamma), \gamma \sim g_2\left(1, \sqrt{\sigma}\right)$$

$$P(Y|\lambda), \lambda \sim g_2\left(\mu, \sqrt{\sigma}\right) = g_2\left(nsc, \sqrt{1/(nc)}\right) = g_1(nc, s)$$

(11).

It means that the gamma distribution goes $g_1(c, s) \Rightarrow g_1(nc, s)$, or $g_2\left(sc, \sqrt{1/c}\right) \Rightarrow g_2\left(nsc, \sqrt{1/(nc)}\right)$ before and after integral across samples.