$$\pi_{ij} = \frac{X_{ij}}{M_j} \qquad \text{random proportion of gene } i \text{ mRNA in observation } j$$

$$E[\pi_{ij}] = E\left[\frac{X_{ij}}{M_j}\right]$$

$$= E\left[E\left[\frac{X_{ij}}{M_j} \mid M_j\right]\right]$$

$$= E\left\{E\left[\frac{a_i M_j}{M_j} \mid M_j\right]\right\}$$

$$= E[a_i] = a_i$$

$$\text{Var}(\pi_{ij}) = E\left[\text{Var}(\pi_{ij} \mid M_j)\right] + \underbrace{\text{Var}\left(\overbrace{E[\pi_{ij} \mid M_j]}^{=a_i}\right)}_{\underset{0}{d}}$$

$$\begin{array}{l} \text{Var}(cX) \\ = c^2 \text{Var}(X) \end{array}$$

$$\text{Var}(\pi_{ij} \mid M_j) = \text{Var}\left(\frac{X_{ij}}{M_j} \mid M_j\right)$$

$$= \frac{1}{M_j^2} \text{Var}(X_{ij} \mid M_j)$$

$$\approx \frac{1}{M_j^2} a_i M_j = \frac{a_i}{M_j}$$

$$\Rightarrow \text{Var}(\pi_{ij}) \approx E\left[\frac{a_i}{M_j}\right] = a_i E\left[\frac{1}{M_j}\right]$$

$$\approx \text{"biological variance"}$$

## Step 2.

This is assumed (here) to be completely random sampling.

Let $D_j$ be the "read depth" of observation $j$, which is the total number of reads from observation $j$.

Note: We observe $D_j$, say $d_j$.

Reminder: $Y_{ij}$ RNA-Seq read counts for gene $i$, obs. $j$.

$$Y_{ij} | \pi_{ij}, d_j \sim \text{Binomial}(d_j, \pi_{ij})$$

$$\sim \text{Poisson}(d_j \pi_{ij})$$

$$d_j \text{ large}, \quad \pi_{ij} \text{ small}$$

$$E[Y_{ij}] = E\left[E[Y_{ij} | \pi_{ij}, d_j]\right]$$

$$= E[\pi_{ij} d_j]$$
$$= d_j E[\pi_{ij}]$$
$$= d_j a_i$$

$$Var(Y_{ij}) = E[Var(Y_{ij} | \pi_{ij}, d_j)] + Var\left(E[Y_{ij} | \pi_{ij}, d_j]\right)$$

$$\approx E[\pi_{ij} d_j] + Var(\pi_{ij} d_j)$$

$$= a_i d_j + d_j^2 Var(\pi_{ij})$$

$$= a_i d_j + d_j^2 a_i E\left[\frac{1}{m_j}\right]$$

$$> E[Y_{ij}] \implies \text{over-dispersed Poisson}$$

$$\text{Estimate } \hat{\pi}_{ij} = \frac{Y_{ij}}{d_j} \quad , \quad E[\hat{\pi}_{ij}] = a_i$$

$$Var(\hat{\pi}_{ij}) = \frac{1}{d_j^2} Var(Y_{ij})$$

$$= \frac{a_i}{d_j} + Var(\pi_{ij})$$

$$= \frac{a_i}{d_j} + a_i E\left[\frac{1}{m_j}\right]$$

technical variance $\approx \dfrac{a_i}{d_j}$

biological variance $\approx a_i \, E\left[\dfrac{1}{M_j}\right]$

Consider estimate $\hat{a}_i = \dfrac{\sum\limits_{j=1}^{n} \hat{\pi}_{ij}}{n}$ ,

$E[\hat{a}_i] = a_i$.

$$Var(\hat{a}_i) = Var\left(\dfrac{\sum\limits_{i=1}^{n} \hat{\pi}_{ij}}{n}\right)$$

$$= \dfrac{1}{n^2} Var\left(\sum\limits_{j=1}^{n} \hat{\pi}_{ij}\right)$$

$$= \dfrac{1}{n^2} \sum\limits_{j=1}^{n} Var(\hat{\pi}_{ij})$$

$$\approx \dfrac{a_i}{n^2} \sum\limits_{j=1}^{n} \dfrac{1}{d_j} + \dfrac{\sum\limits_{j=1}^{n} Var(\pi_{ij})}{n^2}$$

Let's assume $M_j$ are i.i.d.

Then $E\left[\dfrac{1}{M_j}\right]$ is the same for all $j$.

Coefficient of Variation, $CV$

$$CV = \frac{\sqrt{Var(\pi_{ij})}}{a_i} \implies \text{biological } CV$$

$$(CV)^2 = \frac{Var(\pi_{ij})}{a_i^2} = \frac{1}{a_i} E\left[\frac{1}{m_j}\right] \equiv \phi_i$$

Let $\mu_{ij} = d_j a_i$ (population mean times obs. read depth)

$$\implies Var(Y_{ij}) = d_j a_i + d_j^2 \, Var(\pi_{ij})$$

$$= d_j a_i + (d_j a_i)^2 \frac{Var(\pi_{ij})}{a_i^2}$$

$$= \mu_{ij} + \mu_{ij}^2 \, \phi_i$$

estimated by "borrowing strength" across genes — with similar $\phi_i$ values

---

Negative Binomial

Bernoulli trials with success $p$

$Y = $ number of failures before $r^{\underline{th}}$ success

$Y \sim Neg Bin (r, p)$

$$Pr(Y=y) = \binom{r+y-1}{y} p^r (1-p)^y$$

$$y = 0, 1, 2, \ldots$$

$$E[Y] = r \frac{(1-p)}{p} \quad Var(Y) = \frac{r(1-p)}{p^2}$$

$$\text{Let } \mu = \frac{r(1-p)}{p} \bigg\}$$

$$Var(Y) = \mu + \mu^2 \boxed{\frac{1}{r}} \underset{\phi}{\smile}$$

RNA-Seq data under the above model is therefore sometimes modeled as a Neg Bin.

$$Y_{ij} \sim \text{Neg Bin} (p_{ij}, r_i)$$

$$\text{where } \mu_{ij} = \frac{r_i(1-p_{ij})}{p_{ij}} \quad \text{and } \phi_i = \frac{1}{r_i}.$$

## Compound Gamma - Poisson Distribution

$$Y_{ij} \mid \lambda_{ij} \sim \text{Poisson}(\lambda_{ij})$$

$$\lambda_{ij} \sim \text{Gamma}(\alpha, \beta)$$

Then $Y_{ij}$ is marginally a Gamma-Poisson r.v.

Neg Bin is a special case of Gamma-Poisson.

Suppose $Y|\lambda \sim \text{Poisson}(\lambda)$

$\lambda \sim \text{Gamma}(\alpha, \beta)$

Gamma pdf

$$f(\lambda; \alpha, \beta) = \frac{\lambda^{\beta-1} e^{-\lambda/\alpha}}{\alpha^\beta \Gamma(\beta)} \qquad \lambda > 0$$

$$E[\lambda] = \alpha\beta \quad, \quad \text{Var}(\lambda) = \alpha^2 \beta$$

Gamma-Poisson:

$$f(y; \alpha, \beta) = \frac{\Gamma(y+\beta)\alpha^y}{\Gamma(\beta)(1+\alpha)^{\beta+y} y!}$$

$$E[Y] = \alpha\beta \qquad \text{Var}(Y) = \alpha\beta + \alpha^2\beta$$

Let $\mu = \alpha\beta$

$$\text{Var}(Y) = \mu + \mu^2 \cdot \frac{1}{\beta}$$

Let's map this back to the $\alpha$-step model.

$$Y_{ij} \mid \lambda_{ij} \sim \text{Poisson}(\lambda_{ij})$$

$$\lambda_{ij} \sim \text{Gamm}(\alpha, \beta) \rightarrow$$ Let's determine $\lambda_{ij}$, what are $\alpha$, and $\beta$

$\lambda_{ij}$ is $\pi_{ij} d_j$, which is random

$$E[\pi_{ij} d_j] = a_i d_j \quad \text{and}$$

$$E[\lambda_{ij}] = \alpha \beta \cdot \text{so}$$

$$\mu_{ij} = a_i d_j = \alpha \beta$$

$$\text{Var}(Y_{ij}) = \mu_{ij} + \mu_{ij}^2 \phi_i,$$

$$\text{so} \quad \frac{1}{\beta} = \phi_i = \frac{1}{a_i} E\left[\frac{1}{M_j}\right].$$

$$\Rightarrow \beta_{ij} = a_i E\left[\frac{1}{M_j}\right]^{-1}$$

$$\Rightarrow \alpha_{ij} = a_i d_j \cdot \frac{1}{a_i} E\left[\frac{1}{M_j}\right]$$

$$= d_j \cdot E\left[\frac{1}{M_j}\right]$$

$$\Rightarrow \lambda_{ij} \sim \text{Gamma}(\alpha_{ij}, \beta_{ij}).$$

$$Y \mid \lambda \sim \text{Poisson}(\lambda) \quad , \quad \lambda > 0$$

$\lambda$ is r.v. with $E[\lambda]$ and $\text{Var}(\lambda)$

$$\text{Var}(Y) = E[\text{Var}(Y \mid \lambda)] + \text{Var}(E[Y \mid \lambda])$$

$$= E[\lambda] + \text{Var}(\lambda)$$

Since $\lambda > 0$, $E[\lambda]$ will "appear"

in $\text{Var}(\lambda) \implies$ mean-variance relationship

## Sums of Random Variables

If $X$ is a rv and $a, b$ are constants,

then $E[a + bX] = a + bE[X]$ and

$\text{Var}(a + bX) = b^2 \text{Var}(X)$.

Let $X_1, X_2, \ldots, X_n$ be $n$ rv's. then

$$E\left[\sum_{i=1}^{n} X_i\right] = \sum_{i=1}^{n} E[X_i]$$

$$\text{Var}\left(\sum_{i=1}^{n} X_i\right) = \sum_{i=1}^{n} \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j)$$

$$= \sum_{i=1}^{n} \sum_{j=1}^{n} \text{Cov}(X_i, X_j)$$

When $X_1, X_2, \ldots, X_n$ are independent,

then $\text{Cov}(X_i, X_j) = 0$ for $i \neq j$

So $\text{Var}(\sum X_i) = \sum \text{Var}(X_i)$.

Let $\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$. Suppose $X_1, X_2, \ldots,$

$X_n$ are independent.

$$E[\overline{X}_n] = E\left[\frac{1}{n} \sum_{i=1}^{n} X_i\right]$$

$$= \frac{1}{n} E\left[\sum_{i=1}^{n} X_i\right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} E(X_i)$$

when $E[X_1] = E[X_2] = \cdots = E[X_n] = \theta$

then $E[\overline{X}_n] = \theta$.

$$\text{Var}(\overline{X}_n) = \frac{1}{n^2} \sum \text{Var}(X_i)$$

If $\text{Var}(X_1) = \text{Var}(X_2) = \cdots = \text{Var}(X_n) = \tau^2$

$$\text{Var}(\overline{X}_n) = \frac{\tau^2}{n}.$$

## Convergence of rv's

Let $Z_1, Z_2, \ldots$ be a sequence of rv's.

Examples: $Z_n = \bar{X}_n$.

$$Z_n \sim \text{Binomial}(n, p)$$

① Convergence in Distribution:

$\{Z_n\}$ converges in distribution to rv $W$

$$Z_n \xrightarrow{D} W \quad \text{as } n \to \infty \text{ if}$$

$$F_{Z_n}(y) = Pr(Z_n \leq y) \longrightarrow Pr(W \leq y) = F_W(y)$$

for all $y \in \mathbb{R}$, as $n \to \infty$.

② Convergence in Probability:

$$Z_n \xrightarrow{P} W \quad \text{as } n \to \infty \text{ if}$$

$$Pr(|Z_n - W| \leq \varepsilon) \longrightarrow 1 \quad \text{as}$$

$$n \to \infty, \quad \text{for } \varepsilon > 0.$$

{ If $\theta$ is a fixed number, then we can have $Z_n \xrightarrow{P} \theta$.

③ Almost sure convergence

$\{Z_n\}$ converges "almost surely" (a.s.) or "with probability 1" to $W$

$$Z_n \xrightarrow{a.s.} W \quad \text{if}$$

$$\Pr\left(\{w: |Z_n(w) - W(w)| \xrightarrow{n \to \infty} 0\}\right) = 1$$

## Strong Law of Large Numbers

Suppose $X_1, X_2, \ldots$ are iid rv's with population mean $E[X_i] = \mu$ where $E[|X_i|] < \infty$. Then

$$\bar{X}_n \xrightarrow{a.s.} \mu, \quad \text{as } n \to \infty.$$

## Central Limit Theorem

Suppose $X_1, X_2, \ldots$ are iid rv's with $E[X_i] = \mu$ and $Var(X_i) = \sigma^2$.

Then as $n \to \infty$,

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} \text{Normal}(0, \sigma^2)$$

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} \text{Normal}(0, 1)$$

Note:

$$\text{Var}(\bar{X}_n - \mu) = \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

$$\text{Var}\left(\frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}}\right) = \frac{1}{\sigma^2/n} \text{Var}(\bar{X}_n) = 1$$

$$\frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} = \sqrt{n}\left(\frac{\bar{X}_n - \mu}{\sigma}\right) \xrightarrow{D} N(0, 1).$$

## Normal rv's — some useful facts

$$X_1, X_2, \ldots, X_n \stackrel{iid}{\sim} \text{Normal}(\mu, \sigma^2)$$

$$E[\bar{X}_n] = \mu, \quad \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

$$\bar{X}_n \sim \text{Normal}(\mu, \sigma^2/n)$$

because

$$X_1 + X_2 + \cdots + X_n \sim Normal(n\mu, n\sigma^2)$$

$$aX_1 + b \sim Normal(a\mu + b, a^2\sigma^2)$$