

# Primer on Negative Binomial in the context of RNAseq analysis

Albert Lee  
Columbia University  
Rabadan Lab

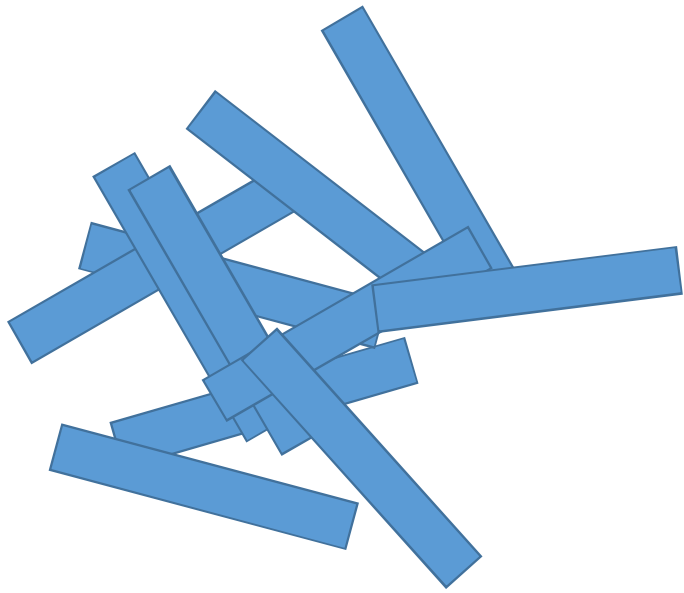
# Outline

- Overview on analyzing count data
- Binomial  $\rightarrow$  Poisson  $\rightarrow$  Negative binomial (NB)
- Motivation for NB in RNA-seq analysis
- Derivation of NB from the marginalization of Gamma-Poisson mixture
- Mean / Variance relationship of Negative Binomial
- Why this framework works well



$n$  flips

$\theta$  chance to see a head



$n$  cDNA fragments  
(reads)

Gene  $g$

$\theta_g$  chance to map to the gene  $g$

Let's define

$X =$  the number of reads mapped to gene  $g$  out of  $n$  total reads,  
where each read has the probability  $\theta_g$   
then

Let's define

$X =$  the number of reads mapped to gene  $g$  out of  $n$  total reads,  
where each read has the probability  $\theta_g$   
then

$$X \sim \text{Bin}(n, \theta_g)$$

$$P(X = x) = \binom{n}{x} \theta_g^x (1 - \theta_g)^{n-x}$$

But what if you don't know  $n$ ?  
(total number of reads)

When  $n \rightarrow \infty$

$$\textit{Bin}(n, \theta) \rightarrow \textit{Pois}(\lambda)$$

Where  $\lambda = n\theta$

$$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$$



# Proof

$$P(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$$

Define  $\lambda = n\theta$

$$= \frac{n!}{(n-x)! x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$\Rightarrow \theta = \frac{\lambda}{n}$$

$$= \frac{n(n-1) \dots (n-x+1)(n-x)!}{(n-x)! x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$= \frac{n(n-1) \dots (n-x+1)}{x!} \left(\frac{\lambda^x}{n^x}\right) \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$= \frac{n^x \left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{x}{n} + \frac{1}{n}\right)}{x!} \left(\frac{\lambda^x}{n^x}\right) \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$= \frac{n^x \left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{x}{n} + \frac{1}{n}\right)}{x!} \left(\frac{\lambda^x}{n^x}\right) \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

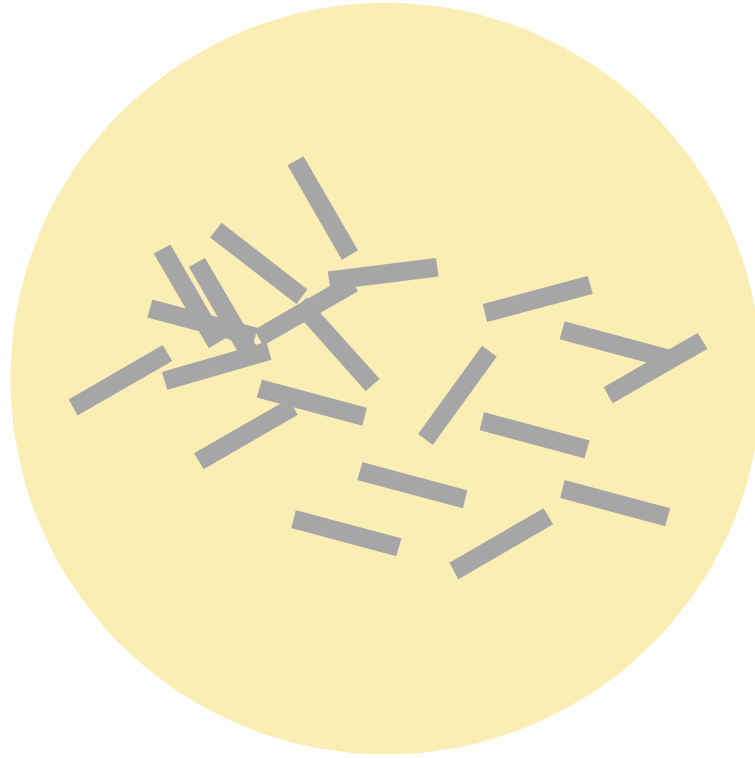
$$= \frac{\left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{x}{n} + \frac{1}{n}\right)}{x!} \lambda^x \left(1 - \frac{\lambda}{n}\right)^{n-x}$$

$$= \frac{\left(1 - \frac{1}{n}\right) \cdots \left(1 - \frac{x}{n} + \frac{1}{n}\right)}{x!} \lambda^x \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}$$

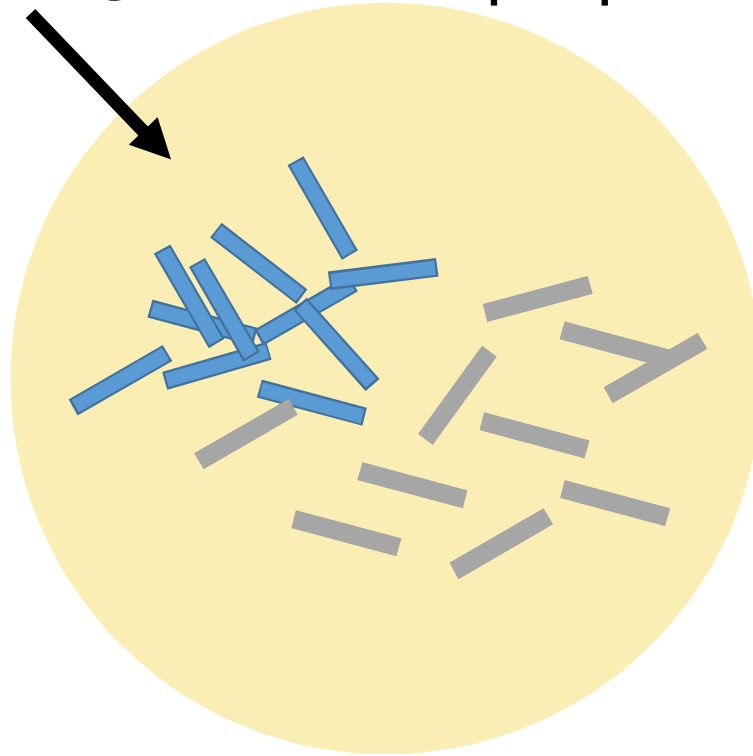
$$n \rightarrow \infty$$

$$= \frac{1}{x!} \lambda^x e^{-\lambda}$$

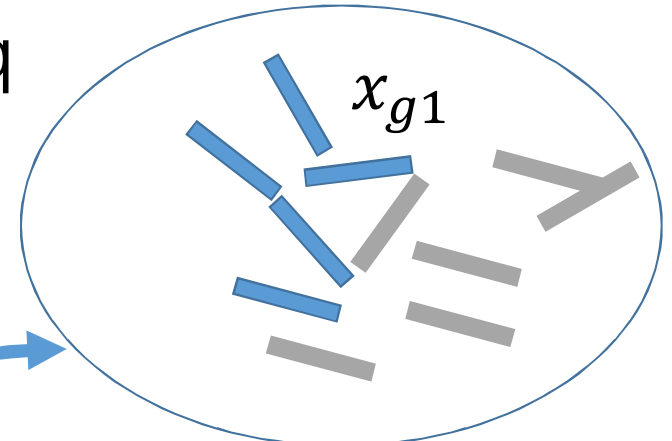
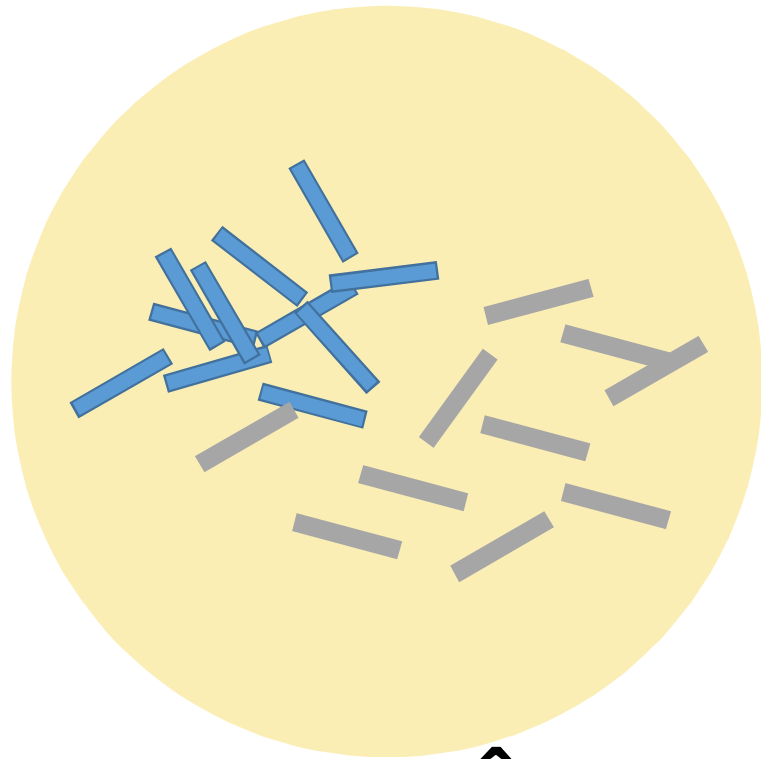
$$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n$$



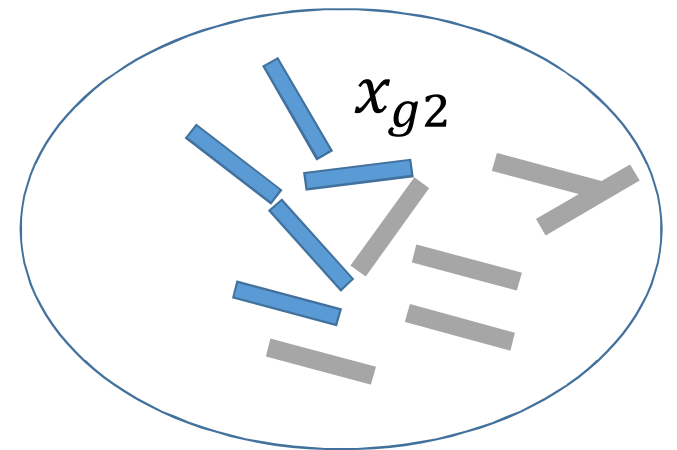
Gene  $g$  With the proportion  $\theta_g$  or  $\lambda_g$



Sampling via RNA-seq



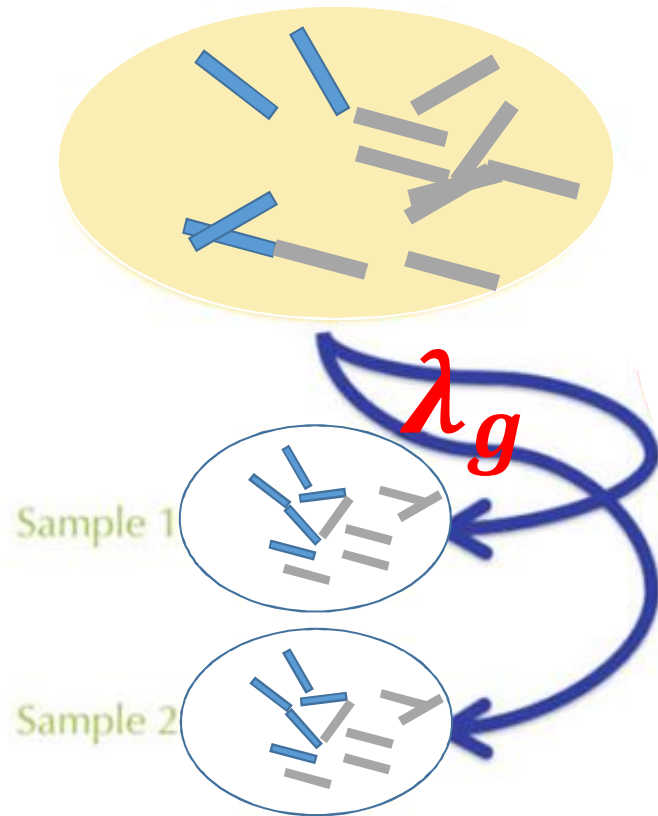
Replicate 1



Replicate 2

$$\lambda_g \approx \hat{\lambda}_g = \frac{\sum x_g}{\# \text{ replicate}}$$

cDNA in the sample  $i$  mappable to gene  $g$



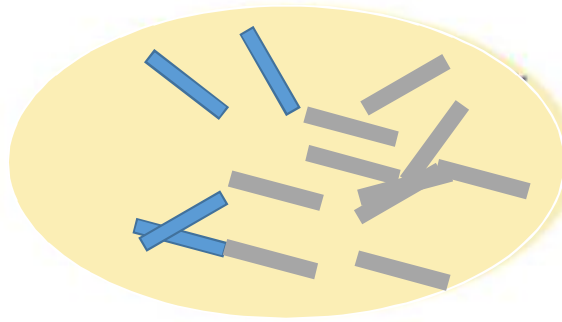
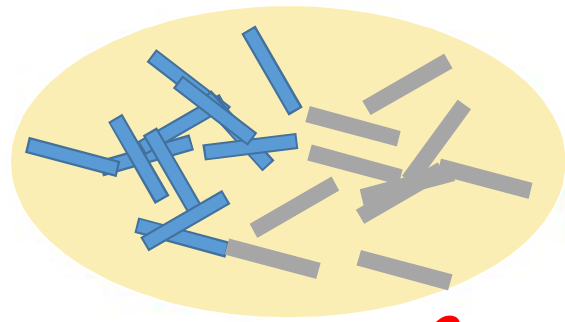
But this only reflects when  $\theta_g$  (or  $\lambda_g$ ) is constant

=> Technical replicates

But..

Fragments in the sample  $i$   
mappable to gene  $g$

Fragments in the sample  $j$   
mappable to gene  $g$

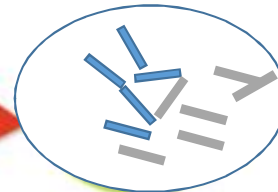


$\lambda_{gi}$

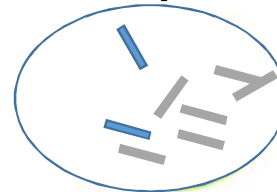
$\lambda_{gj}$

When samples come from  
Different biological replicates,

$$\lambda_{gi} \neq \lambda_{gj}$$



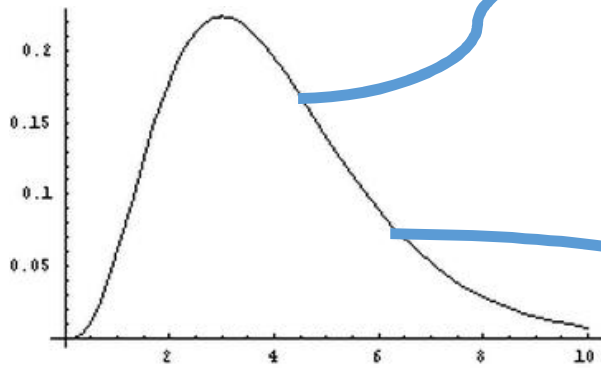
Sample  $i$



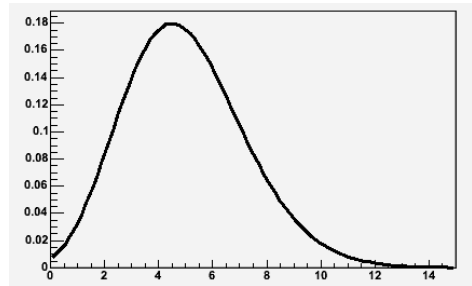
Sample  $j$

# What do you do?

What distribution should we use?  
For this?

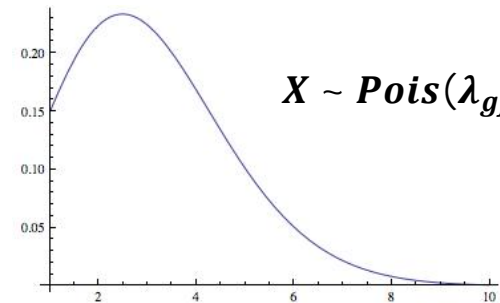


$\lambda_{gi}$



$X \sim \text{Pois}(\lambda_{gi})$

$\lambda_{gj}$



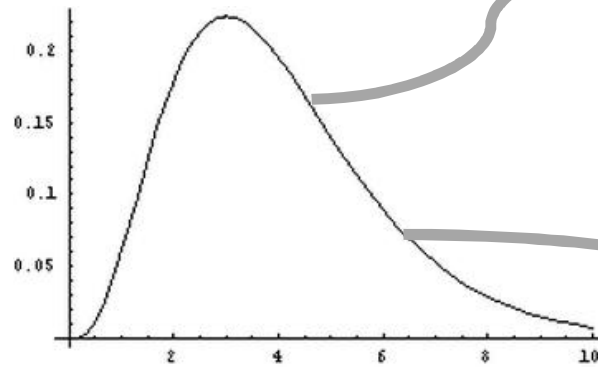
$X \sim \text{Pois}(\lambda_{gj})$

Make  $\lambda$  a random variable!



# We use gamma distribution

What distribution should we use?  
For this?



$\lambda_{gi}$

$\lambda_{gj}$

$$P_{gamma}(X = x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$$

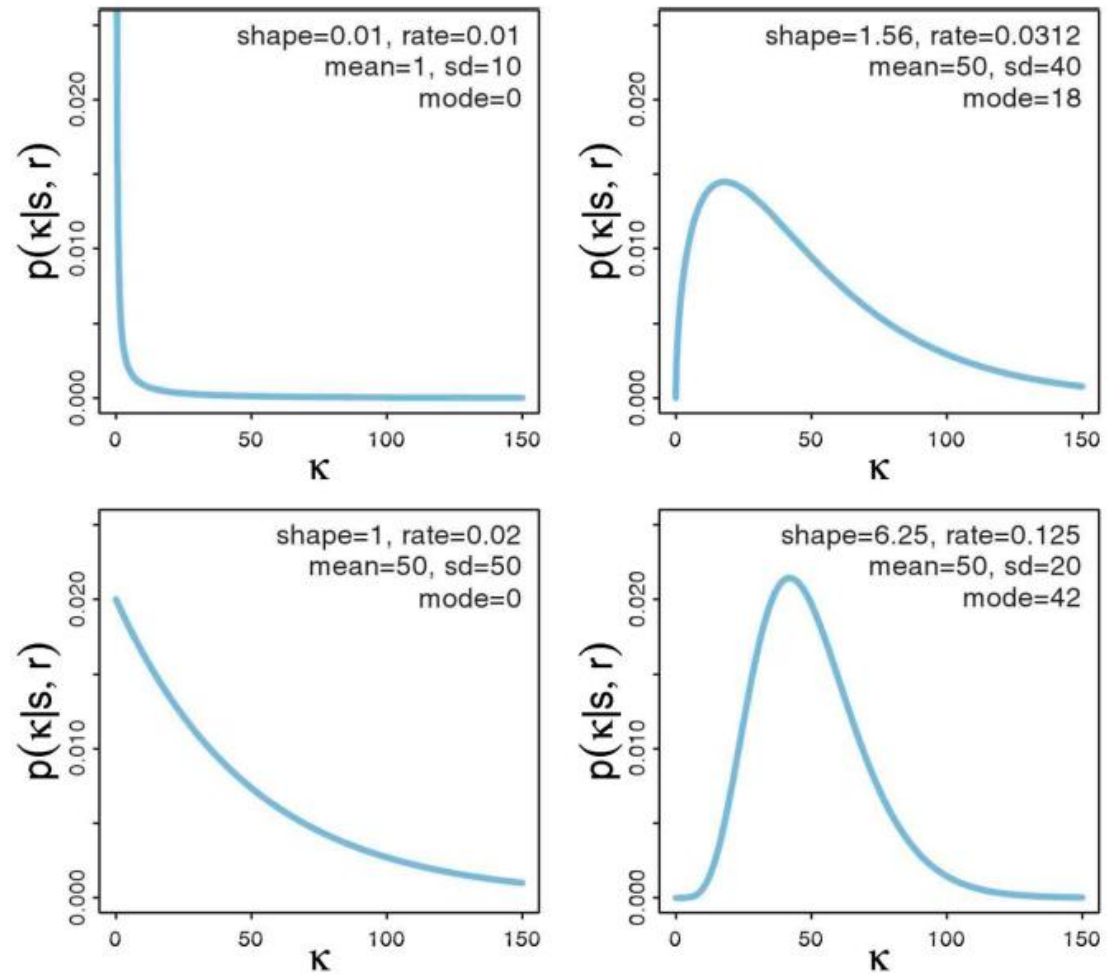
# Why gamma?

- It's mathematically convenient

$$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$$

Shape parameter  $\alpha$

Rate parameter  $\beta$

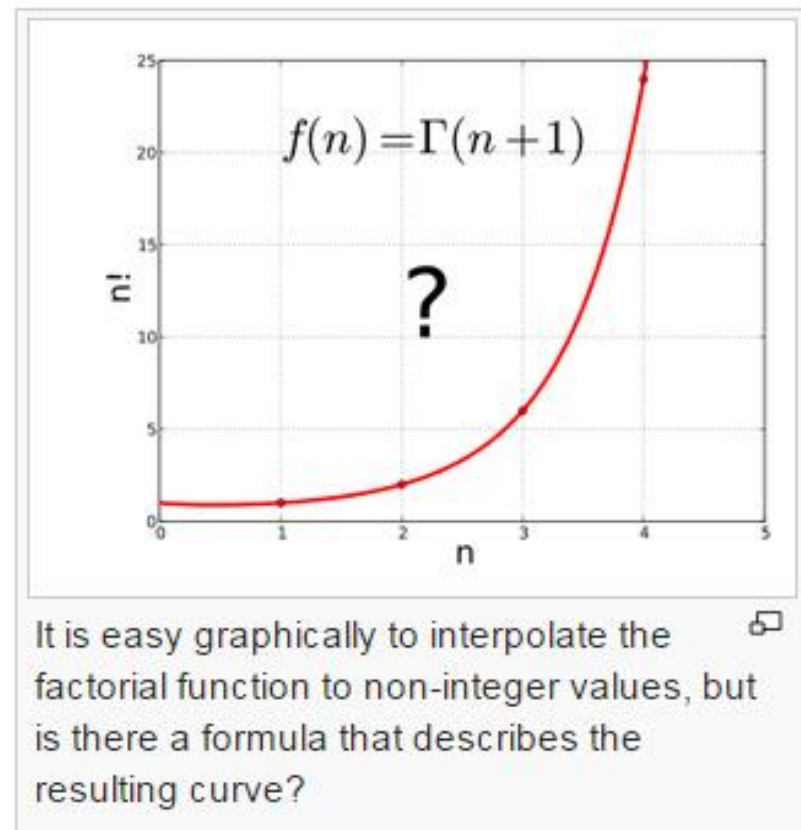


**Figure 9.8** Examples of the gamma distribution. The vertical axis is  $p(\kappa|s, r)$  where  $s$  is the shape and  $r$  is the rate, whose values are annotated in each panel. From Doing Bayesian Analysis 2<sup>nd</sup> ed

What is Gamma function?  $\Gamma(t)$  in  $\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$

$$\Gamma(t) = (t-1)! \quad \text{if } t \in \mathbb{Z} \text{ \& } t-1 \geq 0$$

but what if  $t \in \mathbb{R}$  \&  $t-1 \geq 0$



# What is Gamma function? $\Gamma(t)$

$$\Gamma(t) = \int_0^{\infty} x^{t-1} e^{-x} dx$$

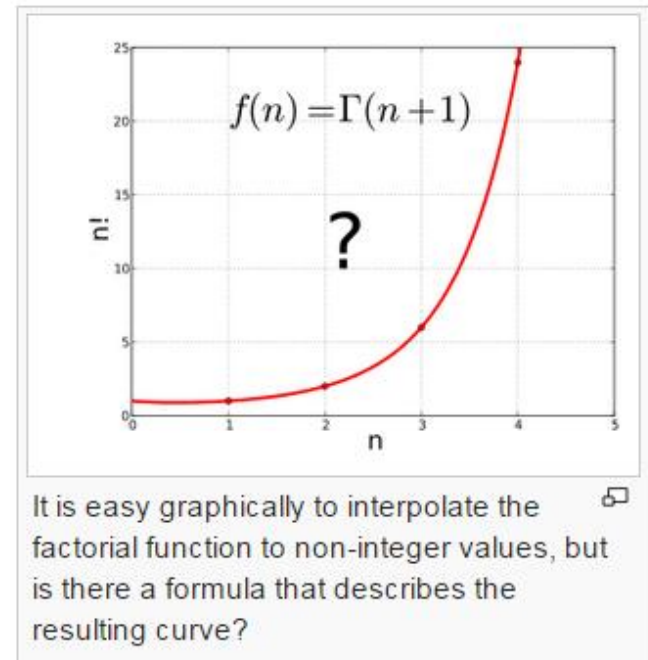
$$\frac{\Gamma(t)}{a^t} = \int_0^{\infty} x^{t-1} e^{-ax} dx$$

$$\int_0^{\infty} x^{t-1} e^{-(ax)} dx = a^{t-1} a^{1-t} \int_0^{\infty} x^{t-1} e^{-(ax)} dx$$

$$= a^{1-t} \int_0^{\infty} (ax)^{t-1} e^{-(ax)} dx$$

$$= a^{-t} \int_0^{\infty} u^{t-1} e^{-u} du = a^{-t} \Gamma(t)$$

let  $u = ax$   
 $du = a dx$

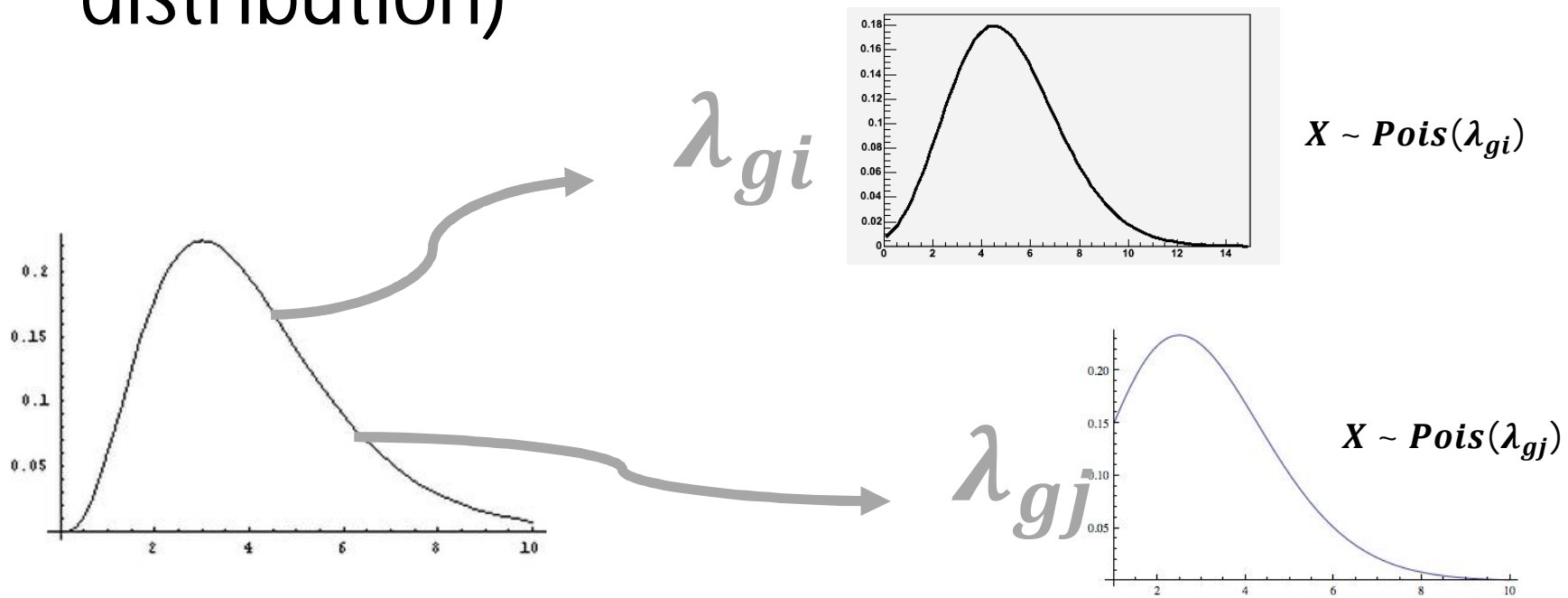


Why Gamma distribution is valid prob. Dist?

$$P_{\text{gamma}}(X = x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x)$$

$$\begin{aligned} \int P_{\text{gamma}}(X = x) &= \int \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} \exp(-\beta x) dx \\ &= \frac{\beta^\alpha}{\Gamma(\alpha)} \int x^{\alpha-1} \exp(-\beta x) dx \\ &= \frac{\beta^\alpha}{\Gamma(\alpha)} \beta^{-\alpha} \Gamma(\alpha) = 1 \end{aligned}$$

# Hierarchical modeling (or Mixture distribution)



$$\Lambda \sim \Gamma(a, \beta)$$
$$X \mid \Lambda = \lambda \sim \text{Pois}(\lambda)$$

# Gamma-Poisson mixture

$$\begin{aligned}\Lambda &\sim \Gamma(\alpha, \beta) && \rightarrow P(\lambda) \\ (X \mid \Lambda = \lambda) &\sim \text{Pois}(\lambda) && \rightarrow P(x \mid \lambda)\end{aligned}$$

$$P(X = x) = ?$$

$$P(X = x) = \int_0^{\infty} P(\lambda, x) d\lambda \quad \text{marginalization}$$

$$= \int_0^{\infty} P(x \mid \lambda) P(\lambda) d\lambda \quad \text{Chain rule}$$

$$= \int_0^{\infty} \left( \frac{\lambda^x e^{-\lambda}}{x!} \right) \left( \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{(-\beta\lambda)} \right) d\lambda$$

$$= \int_0^{\infty} \left( \frac{\lambda^x e^{-\lambda}}{x!} \right) \left( \frac{\beta^\alpha}{\Gamma(a)} \lambda^{a-1} e^{-\beta\lambda} \right) d\lambda$$

$$= \frac{1}{x!} \frac{\beta^\alpha}{\Gamma(a)} \int_0^{\infty} \lambda^x e^{-\lambda} \lambda^{a-1} e^{-\beta\lambda} d\lambda$$

$$= \frac{1}{x!} \frac{\beta^\alpha}{\Gamma(a)} \int_0^{\infty} \lambda^{(x+a)-1} e^{-(\beta+1)\lambda} d\lambda$$

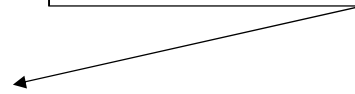
$$= \frac{1}{x!} \frac{\beta^\alpha}{\Gamma(a)} \left( \frac{(\beta+1)^{x+a}}{\Gamma(x+a)} \right)^{-1} \left( \frac{(\beta+1)^{x+a}}{\Gamma(x+a)} \right) \int_0^{\infty} \lambda^{(x+a)-1} e^{-(\beta+1)\lambda} d\lambda$$

$$= \frac{1}{x!} \frac{\beta^\alpha}{\Gamma(a)} \left( \frac{(\beta+1)^{x+a}}{\Gamma(x+a)} \right)^{-1} \left( \frac{(\beta+1)^{x+a}}{\Gamma(x+a)} \right) \int_0^{\infty} \lambda^{(x+a)-1} e^{-(\beta+1)\lambda} d\lambda$$

$$= \frac{1}{x!} \frac{\beta^\alpha}{\Gamma(a)} \left( \frac{(\beta+1)^{x+a}}{\Gamma(x+a)} \right)^{-1} \quad (1)$$

$$= \frac{1}{x!} \frac{\beta^\alpha}{\Gamma(a)} \left( \frac{\Gamma(x+a)}{(\beta+1)^{x+a}} \right) = \frac{\Gamma(x+a)}{x! \Gamma(a)} \left( \frac{1}{\beta+1} \right)^x \left( \frac{\beta}{\beta+1} \right)^a$$

This is negative binomial!





$$\begin{aligned}
&= \frac{\Gamma(x+a)}{x! \Gamma(a)} \left(\frac{1}{\beta+1}\right)^x \left(\frac{\beta}{\beta+1}\right)^a \\
&= \left(\frac{(x+a-1)!}{x! (a-1)!}\right) \left(\frac{1}{\beta+1}\right)^x \left(\frac{\beta}{\beta+1}\right)^a \\
&= \binom{x+a-1}{a-1} \left(\frac{1}{\beta+1}\right)^x \left(\frac{\beta}{\beta+1}\right)^a \\
&= \binom{x+a-1}{a-1} p^x (1-p)^a
\end{aligned}$$

Let's restrict that  $x+a$  is integer, then  
 $\Gamma(x) = (x-1)!$

Because  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$

Define  $p = \frac{1}{\beta+1} \Rightarrow 1-p = \frac{\beta}{\beta+1}$

If  $X$  is the number of success before  $a$  failures, then  
 $X \sim NB(a, p)$

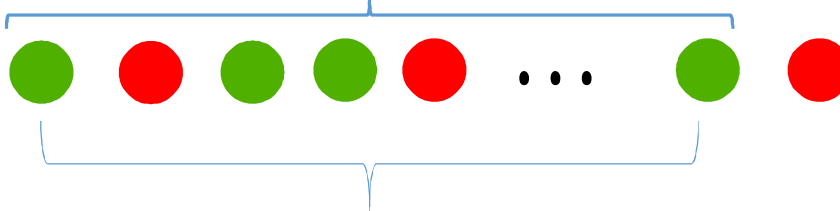
# Why ?

If  $X$  is the number of success before  $a$  failures, then  
 $X \sim NB(a, p)$

$X$  = number of success 

$a$  = number of failures 

But the sequence of the balls before the last step can vary



$(x + a - 1)$  many balls

By definition, failure must happen at the end

$\binom{x + a - 1}{a - 1}$  = total number of possible combination of failures before the last step

$$= \binom{x + a - 1}{a - 1} p^x (1 - p)^{a-1} (1 - p)$$

$$= \binom{x + a - 1}{a - 1} p^x (1 - p)^a$$

We know that for Poisson :

$$E[X_{pois}] = \lambda$$

$$Var[X_{pois}] = \lambda$$

But for NB, what's

$$E[X_{NB}] = ?$$

$$Var[X_{NB}] = ?$$

# Mean of negative binomial

$$E[X] = \int xP(X)dx$$

$$= \int_0^{\infty} x \frac{\Gamma(x+r)}{x! \Gamma(r)} p^x (1-p)^r dx$$

NASTY !

# Instead solve it via Moment generating function of NB

Moment generating function = magical formula that spits out things we need:  
mean, variance, skewness, kurtosis, ...  
or what's collectively known as  
"moments" (shape of the distribution)

$$M_X(t) = E[e^{tX}]$$

$$M_X'(0) = E[X]$$

Mean

$$M_X''(0) = E[X^2]$$

$$\text{Var}[X] = E[X^2] - E[X]^2$$

Variance

# “Secret of the magic trick”

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} = \frac{x^0}{0!} + \frac{x^1}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$$

$$e^{tX} = \frac{(tX)^0}{0!} + \frac{(tX)^1}{1!} + \frac{(tX)^2}{2!} + \frac{(tX)^3}{3!} + \dots$$

$$M_X(t) = E[e^{tX}] = E\left[\frac{(tX)^0}{0!} + \frac{(tX)^1}{1!} + \frac{(tX)^2}{2!} + \frac{(tX)^3}{3!} + \dots\right]$$

$$= E\left[\frac{(tX)^0}{0!}\right] + E\left[\frac{(tX)^1}{1!}\right] + E\left[\frac{(tX)^2}{2!}\right] + E\left[\frac{(tX)^3}{3!}\right] + \dots$$

Linearity of expectation

$$= 1 + tE[X] + \frac{t^2 E[X^2]}{2!} + \frac{t^3 E[X^3]}{3!} + \dots$$

$E[cX] = cE[X]$

$$M'_X(t) = 0 + E[X] + tE[X^2] + \frac{t^2 E[X^3]}{2!} + \dots \longrightarrow M'_X(0) = E[X]$$

$$M''_X(t) = E[X^2] + tE[X^3] + \dots \longrightarrow M''_X(0) = E[X^2]$$

# Instead solve it via Moment generating function of NB

$$\begin{aligned}
 M_X(t) &= E[e^{tX}] \\
 &= \sum_{x=0}^{\infty} e^{tx} \binom{x+r-1}{x} p^x (1-p)^a \\
 &= (1-p)^a \sum_{x=0}^{\infty} e^{tx} \binom{x+r-1}{x} p^x \\
 &= (1-p)^a \sum_{x=0}^{\infty} \binom{x+r-1}{x} (pe^t)^x \\
 &= \frac{(1-p)^a}{(1-pe^t)^a}
 \end{aligned}$$

$$M_{X_{NB}}(t) = \left( \frac{1-p}{1-pe^t} \right)^r$$

$$\sum_{x=0}^{\infty} \binom{x+a-1}{x} (pe^t)^x = \frac{1}{(1-pe^t)^a}$$

sum\_{(x=0)^{\infty}} e^{tx} choose(x+r-1,x) \* p^x \* (1-p)^r

Examples Random

Input interpretation:

$$\sum_{x=0}^{\infty} e^{tx} \binom{x+r-1}{x} p^x (1-p)^r$$

$\binom{n}{m}$  is the binomial coefficient

Infinite sum:

$$\sum_{x=0}^{\infty} (1-p)^r p^x e^{tx} \binom{x+r-1}{x} = (1-p)^r (1-pe^t)^{-r} \approx (1-p)^r (1-p \cdot 2.71828^t)^{-r}$$

when  $|p| < e^{-\text{Re}(t)}$

$\binom{n}{m}$  is the binomial coefficient  
 $|z|$  is the absolute value of  $z$   
 $\text{Re}(z)$  is the real part of  $z$

# Mean of NB

$$M_X(t) = \left( \frac{1-p}{1-pe^t} \right)^r$$



derivative of  $((1-p)/(1-p \cdot e^t))^a$  with respect to  $t$ ;  $t=0$  ☆

[Examples](#) [Random](#)

Input interpretation:

$$\frac{\partial}{\partial t} \left( \frac{1-p}{1-p e^t} \right)^a \text{ where } t = 0$$

Result:

$$\frac{a p}{1-p}$$

$$M'_X(0) = \frac{ap}{1-p} = E[X]$$



# Variance of NB

WolframAlpha computational knowledge engine

second derivative of  $((1-p)/(1-pe^t))^a$  with respect to  $t$ ;  $t=0$

Input interpretation:  
 $\frac{\partial^2}{\partial t^2} \left( \frac{1-p}{1-pe^t} \right)^a$  where  $t=0$

Result:  
 $\frac{ap(ap+1)}{(p-1)^2}$

$$M_X''(0) = \frac{ap(ap+1)}{(p-1)^2} = E[X^2]$$

$$\text{Var}(X) = E[X^2] - E[X]^2 = \frac{ap(ap+1)}{(p-1)^2} - \left( \frac{ap}{1-p} \right)^2 = \frac{ap}{(1-p)^2}$$

# Denouement

We have  $E[X] = \frac{ap}{1-p}$

$$\text{Var}[X] = \frac{ap}{(1-p)^2}$$

Remember  $p = \frac{1}{\beta + 1}$

$$1 - p = \frac{\beta}{\beta + 1}$$

Then

$$E[X] = \frac{ap}{1-p} = \frac{a \left( \frac{1}{\beta + 1} \right)}{\frac{\beta}{\beta + 1}} = \frac{a}{\beta}$$

$$\begin{aligned} \text{Var}[X] &= \frac{ap}{(1-p)^2} = \frac{a \left( \frac{1}{\beta + 1} \right)}{\left( \frac{\beta}{\beta + 1} \right)^2} = \frac{a}{\beta^2} (\beta + 1) \\ &= \frac{a}{\beta} + \frac{a}{\beta^2} \\ &= \frac{a}{\beta} + \left( \frac{a}{\beta} \right)^2 \left( \frac{1}{a} \right) \end{aligned}$$

Therefore

$$\text{Var}[X] = E[X] + \left( \frac{1}{a} \right) E[X]^2$$

We know that for poisson :

$$E[X_{pois}] = \lambda$$

$$Var[X_{pois}] = \lambda$$

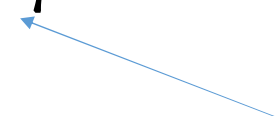
And for NB,

$$E[X_{NB}] = \mu$$

$$Var[X_{NB}] = \mu + c\mu^2$$

  
Dispersion

# Most efforts of most RNAseq packages...

$$E[X_{NB}] = \mu$$
$$Var[X_{NB}] = \mu + c\mu^2$$


Estimate dispersion using small # of samples

**EdgeR** : (old) assume  $c$  is identical for all genes

(new )  $c$  is different for individual genes

Using hierarchical modeling to shrink to  $c$  to a consensus value

**DESeq** : Assume  $c$  and  $\mu$  is related nonlinearly. use local regression to compute  $c$

# Why is this nice???

- Negative binomial is exponential family! So you can model complicated stuff using generalized linear model.

$$E[X_{NB}] = \mu$$

$$\text{Assume } \mu = e^{X^T \beta}$$

Estimate  $\beta$

$X$  = design matrix

$\beta$  = vector of coefficients

- Empirically, packages that use NB as an underlying model are shown to work better than those that don't

THE END