# Variance-stabilizing transformation for DESeq

### For parametrized dispersion fit

This file describes the variance stabilizing transformation (VST) used by DESeq when parametric dispersion estimation is used.
This is a *Mathematica* notebook. The file *vst.pdf* is produced from *vst.nb*.

When using *estimateDispersions* with *fitType="parametric"*, we parametrize the relation between mean $\mu$ and dispersion $\alpha$ with two constants $a_0$ and $a_1$ as follows:

$$\alpha = a_0 + a_1 / \mu$$

$$a_0 + \frac{a_1}{\mu}$$

In the package, $a_0$ is called the *asymptotic dispersion* and $a_1$ the *extra-Poisson factor*.

The variance is hence

$$v = \mu + \alpha \mu^2 \text{ // Expand}$$

$$\mu + \mu^2 a_0 + \mu a_1$$

A variance stabilizing transformation (VST) is a transformation $u$, such that, if $X$ is a random variable with variance-mean relation $v$, i.e., $\mathrm{Var}(X) = v(\mathrm{E}(X))$, then $u(X)$ has stabilized variance, i.e., is homoskedastic.

A VST $u$ can be derived from a variance-mean relation $v$ by $u(x) = \int^x \frac{d\mu}{\sqrt{v(\mu)}}$.

Hence, we can get a general VST with

$$u_0 = \text{Integrate}\left[ \frac{1}{\sqrt{v}}, \{\mu, 0, x\}, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}\right]$$

$$\frac{\text{Log}\left[\frac{1+2 x a_0 + a_1 + 2\sqrt{x a_0 (1+x a_0 + a_1)}}{1+a_1}\right]}{\sqrt{a_0}}$$

If $u_0$ is a VST, then so is $u(x) = \eta u_0(x) + \xi$. Hence, this here is a VST, too:

$$u = \eta u_0 + \xi$$

$$\xi + \frac{\eta \text{ Log}\left[\frac{1+2 x a_0 + a_1 + 2\sqrt{x a_0 (1+x a_0 + a_1)}}{1+a_1}\right]}{\sqrt{a_0}}$$

We will now choose the parameters $\eta$ and $\xi$ such that our VST behaves like $\log_2$ for large values. Let us first look at the asymptotic ratio of the two transformations:

$$\text{Limit}[u / \text{Log}[2, x], x \rightarrow \infty, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}]$$

$$\frac{\eta \text{ Log}[2]}{\sqrt{a_0}}$$

Hence, if we set $\eta$ as follows, both tranformations have asymptotically the ratio 1.

$$\eta = \frac{\sqrt{a_0}}{\text{Log}[2]}$$

$$\frac{\sqrt{a_0}}{\text{Log}[2]}$$

We also want the difference to vanish for large values:

```
Limit[u - Log[2, x], x → ∞, Assumptions → {a₀ > 0, a₁ > 0, x > 0}]
```

$$\xi + \frac{\text{Log}\left[\frac{4\,a_0}{1+a_1}\right]}{\text{Log}[2]}$$

So, we set

$$\xi = -\frac{\text{Log}\left[\frac{4\,a_0}{1+a_1}\right]}{\text{Log}[2]}$$

$$-\frac{\text{Log}\left[\frac{4\,a_0}{1+a_1}\right]}{\text{Log}[2]}$$

Check that both limits are now correct:

```
Limit[u / Log[2, x], x → ∞, Assumptions → {a₀ > 0, a₁ > 0, x > 0}]
```

1

```
Limit[u - Log[2, x], x → ∞, Assumptions → {a₀ > 0, a₁ > 0, x > 0}]
```
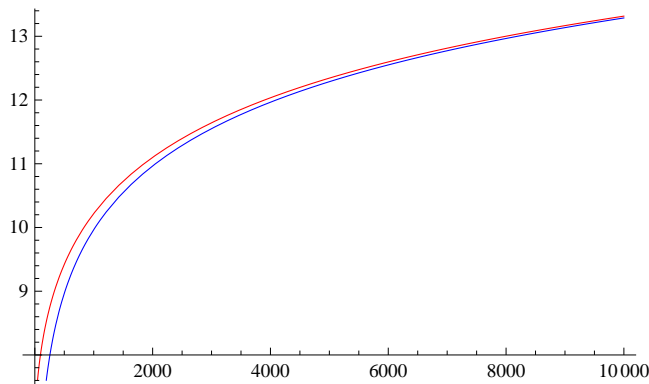
0

Hence, we arrive at this VST:

```
FullSimplify[u, Assumptions -> {a₀ > 0, a₁ > 0, x > 0}]
```

$$\frac{\text{Log}\left[\frac{1+2\,x\,a_0+a_1+2\,\sqrt{x\,a_0\,(1+x\,a_0+a_1)}}{4\,a_0}\right]}{\text{Log}[2]}$$
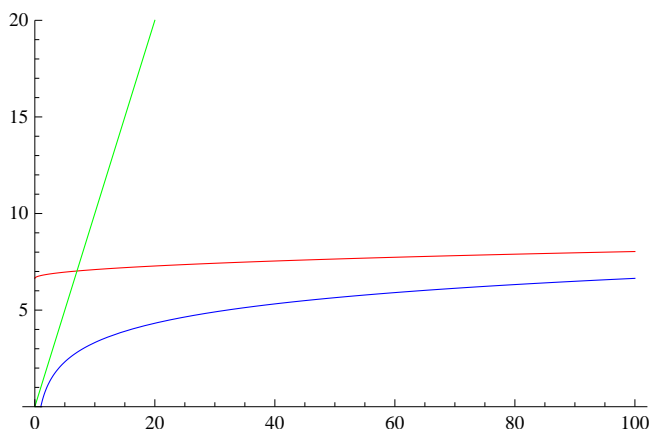
This VST (red) now behaves asymptotically as $\log_2$ (blue), shown here for typical values for $a_0$ and $a_1$.

```
Plot[ {u /. {a₀ → .01, a₁ -> 3}, Log[2, x]}, {x, 0, 10 000}, PlotStyle → {Red, Blue}]
```



For small values, however, the VST (red) compresses the dynamics much more dramatically than the logarithm (blue) and the identity (green). This reflects that the strong Poisson noise makes differences uninformative for small values.

```
Plot[ {u /. {a₀ → .01, a₁ -> 3}, Log[2, x], x},
  {x, 0, 100}, PlotStyle → {Red, Blue, Green}, PlotRange → {0, 20}]
```



A template for the R code in the function:

```
CForm[FullSimplify[u, Assumptions -> {a₀ > 0, a₁ > 0, x > 0}]] /.
  {a₀ → asymptDisp, a₁ → extraPois, x → q}
Log((1 + extraPois + 2*asymptDisp*q +
      2*Sqrt(asymptDisp*q*(1 + extraPois + asymptDisp*q)))/
      (4.*asymptDisp))/Log(2)
```

## For local dispersion fit

In case of a local dispersion fit, the variance-stabilizing transformation $u(x) = \int^x \frac{d\mu}{\sqrt{v(\mu)}}$ is obtained by numerical integration of the fitted mean-dispersion relation $v(\mu)$ (by adding up along a asinh-spaced grid and a fitting a spline). Then, the scaling parameters $\eta$ and $\xi$ (see above) are chosen such that the VST is equal to $\log_2$ for two large normalized count values (for which the 95- and the 99.9-percentile of the sample-averaged normalized count values are used.)