

# Variance-stabilizing transformation for parametrized dispersion

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:

$$\text{In}[1]:= \alpha = a_0 + a_1 / \mu$$

$$\text{Out}[1]= a_0 + \frac{a_1}{\mu}$$

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

The variance is hence

$$\text{In}[2]:= v = \mu + \alpha \mu^2 // \text{Expand}$$

$$\text{Out}[2]= \mu + \mu^2 a_0 + \mu a_1$$

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

$$\text{In}[3]:= u_0 = \text{Integrate}\left[\frac{1}{\sqrt{v}}, \{\mu, 0, x\}, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}\right]$$

$$\text{Out}[3]= \frac{\text{Log}\left[\frac{1+2 x a_0+a_1+2 \sqrt{x a_0 (1+x a_0+a_1)}}{1+a_1}\right]}{\sqrt{a_0}}$$

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

$$\text{In}[4]:= u = \eta u_0 + \xi$$

$$\text{Out}[4]= \xi + \frac{\eta \text{Log}\left[\frac{1+2 x a_0+a_1+2 \sqrt{x a_0 (1+x a_0+a_1)}}{1+a_1}\right]}{\sqrt{a_0}}$$

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:

$$\text{In}[5]:= \text{Limit}[u / \text{Log}[2, x], x \rightarrow \infty, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}]$$

$$\text{Out}[5]= \frac{\eta \text{Log}[2]}{\sqrt{a_0}}$$

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

$$\text{In[6]:= } \eta = \frac{\sqrt{a_0}}{\text{Log}[2]}$$

$$\text{Out[6]= } \frac{\sqrt{a_0}}{\text{Log}[2]}$$

We also want the difference to vanish for large values:

$$\text{In[7]:= } \text{Limit}[u - \text{Log}[2, x], x \rightarrow \infty, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}]$$

$$\text{Out[7]= } \xi + \frac{\text{Log}\left[\frac{4 a_0}{1+a_1}\right]}{\text{Log}[2]}$$

So, we set

$$\text{In[8]:= } \xi = -\frac{\text{Log}\left[\frac{4 a_0}{1+a_1}\right]}{\text{Log}[2]}$$

$$\text{Out[8]= } -\frac{\text{Log}\left[\frac{4 a_0}{1+a_1}\right]}{\text{Log}[2]}$$

Check that both limits are now correct:

$$\text{In[9]:= } \text{Limit}[u / \text{Log}[2, x], x \rightarrow \infty, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}]$$

$$\text{Out[9]= } 1$$

$$\text{In[10]:= } \text{Limit}[u - \text{Log}[2, x], x \rightarrow \infty, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}]$$

$$\text{Out[10]= } 0$$

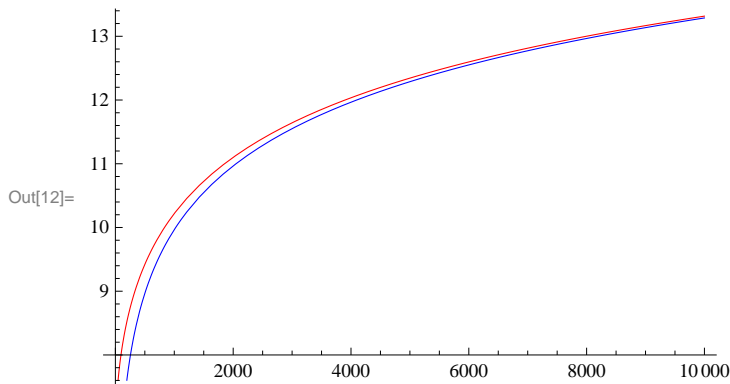
Hence, we arrive at this VST:

$$\text{In[11]:= } \text{FullSimplify}[u, \text{Assumptions} \rightarrow \{a_0 > 0, a_1 > 0, x > 0\}]$$

$$\text{Out[11]= } \frac{\text{Log}\left[\frac{1+2 x a_0+a_1+2 \sqrt{x a_0 (1+x a_0+a_1)}}{4 a_0}\right]}{\text{Log}[2]}$$

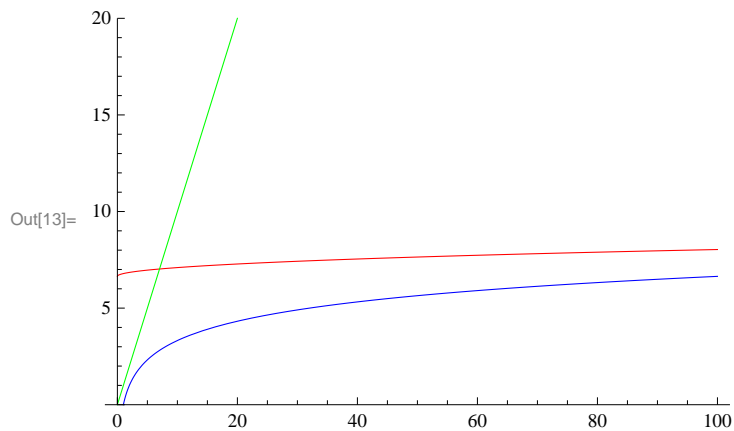
This VST (red) now behaves asymptotically as  $\log_2$  (blue), shown here for typical values for  $a_0$  and  $a_1$ .

$$\text{In[12]:= } \text{Plot}[\{u /. \{a_0 \rightarrow .01, a_1 \rightarrow 3\}, \text{Log}[2, x]\}, \{x, 0, 10000\}, \text{PlotStyle} \rightarrow \{\text{Red}, \text{Blue}\}]$$



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.

```
In[13]:= 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:

```
In[20]:= CForm[FullSimplify[u, Assumptions -> {a0 > 0, a1 > 0, x > 0}]] /.
  {a0 → asympDisp, a1 → extraPois, x → q}
```

```
Out[20]//CForm=
Log((1 + extraPois + 2*asympDisp*q +
  2*Sqrt(asympDisp*q*(1 + extraPois + asympDisp*q)))/
  (4.*asympDisp))/Log(2)
```