

APPENDIX B

The ‘Delta method’ ...

Suppose you have conducted a mark-recapture study over 4 years which yields 3 estimates of apparent annual survival (say, $\hat{\phi}_1$, $\hat{\phi}_2$, and $\hat{\phi}_3$). But, suppose what you are really interested in is the estimate of the product of the three survival values (i.e., the probability of surviving from the beginning of the study to the end of the study)? While it is easy enough to derive an estimate of this product (as $[\hat{\phi}_1 \times \hat{\phi}_2 \times \hat{\phi}_3]$), how do you derive an estimate of the *variance* of the product? In other words, how do you derive an estimate of the variance of a transformation of one or more random variables, where in this case, we transform the three random variables ($\hat{\phi}_1$, $\hat{\phi}_2$ and $\hat{\phi}_3$) by considering their product?

One commonly used approach which is easily implemented, not computer-intensive*, and can be robustly applied in many (but not all) situations is the so-called *Delta method* (also known as the method of propagation of errors). In this appendix, we briefly introduce some of the underlying background theory, and the application of the Delta method.

Applying the ‘Delta method’ effectively assumes you have a good understanding of the underlying concepts and principles. As such, the first few sections of this appendix build up from first principles (i.e., reviews some basic probability and statistical theory). If you’d rather just ‘skip ahead’ to ‘applications’ (i.e., proceed with the mechanics of ‘how to do it’ without fully understanding or caring about ‘why it works’ – and ‘when doesn’t work (and why)’), you might choose to jump ahead to section (B.3).

B.1. Background – mean and variance of random variables

Our interest here is developing a method that will allow us to estimate the variance for functions of random variables. Let’s start by considering the formal approach for deriving these values explicitly, based on the *method of moments*.[†] For continuous random variables, consider a continuous function $f(x)$ on the interval $[-\infty, +\infty]$. The first four moments ($M_0 \rightarrow M_3$) of $f(x)$ can be written as:

$$M_0 = \int_{-\infty}^{+\infty} f(x) dx,$$
$$M_1 = \int_{-\infty}^{+\infty} x f(x) dx,$$

* We briefly discuss some compute-intensive approaches in Addendum (B.1) to this appendix.

[†] In simple terms, a *moment* is a specific quantitative measure, used in both mechanics and statistics, of the shape of a set of points. If the set of points represents a probability density, then the moments relate to measures of shape and location such as mean, variance, skewness, and so forth. We distinguish between *raw* moments, and *central* moments, later on...

$$M_2 = \int_{-\infty}^{+\infty} x^2 f(x) dx,$$

$$M_3 = \int_{-\infty}^{+\infty} x^3 f(x) dx.$$

In the particular case that the function is a probability density (as for a continuous random variable), then $M_0 = 1$ (i.e., the area under the probability density function (pdf) must equal 1).

For example, consider the uniform distribution on the finite interval $\mathcal{U}[a, b]$. A uniform distribution (sometimes also known as a rectangular distribution), is a distribution that has constant probability over the interval. The pdf for a continuous uniform distribution on the finite interval $\mathcal{U}[a, b]$ is:

$$P(x) = \begin{cases} 0 & \text{for } x < a \\ 1/(b - a) & \text{for } a < x < b \\ 0 & \text{for } x > b. \end{cases}$$

Integrating the pdf for $p(x) = 1/(b - a)$:

$$M_0 = \int_a^b p(x) dx$$

$$= \int_a^b \frac{1}{b - a} dx = 1,$$

$$M_1 = \int_a^b xp(x) dx$$

$$= \int_a^b \frac{x}{b - a} dx = \frac{a + b}{2},$$

$$M_2 = \int_a^b x^2 p(x) dx$$

$$= \int_a^b x^2 \frac{1}{b - a} dx = \frac{1}{3}(a^2 + ab + b^2),$$

$$M_3 = \int_a^b x^3 p(x) dx$$

$$= \int_a^b x^3 \frac{1}{b - a} dx = \frac{1}{4}(a^3 + a^2b + ab^2 + b^3).$$

If you look closely, you should be able to deduce that M_1 is the mean of the distribution (i.e., halfway between a and b). But, what about the variance? How do we interpret/use the other moments?

Recall that the variance is defined as the average value of the fundamental quantity [distance from mean]². The squaring of the distance is so the values to either side of the mean don't cancel out. The standard deviation is simply the square-root of the variance.

Given some discrete random variable x_i , with probability p_i , and mean μ , we define the variance as:

$$\text{Var} = \sum (x_i - \mu)^2 p_i.$$

Note we don't have to divide by the number of values of x because the sum of the discrete probability distribution is 1 (i.e., $\sum p_i = 1$).

For a continuous probability distribution, with mean μ , we define the variance as:

$$\text{Var} = \int_a^b (x - \mu)^2 p(x) dx.$$

Next we expand and do a bit of simple algebraic rearrangement of this expression:

$$\begin{aligned} \text{Var} &= \int_a^b (x - \mu)^2 p(x) dx \\ &= \int_a^b (x^2 - 2\mu x + \mu^2) p(x) dx \\ &= \int_a^b x^2 p(x) dx - \int_a^b 2\mu x p(x) dx + \int_a^b \mu^2 p(x) dx \\ &= \int_a^b x^2 p(x) dx - 2\mu \int_a^b x p(x) dx + \mu^2 \int_a^b p(x) dx. \end{aligned}$$

If we look closely at the last line, we see that in fact each of the 3 terms are functions of the different moments of the uniform distribution (derived on the preceding page).

Thus we can write:

$$\begin{aligned} \text{Var} &= \int_a^b (x - \mu)^2 p(x) dx \\ &= \int_a^b x^2 p(x) dx - 2\mu \int_a^b x p(x) dx + \mu^2 \int_a^b p(x) dx \\ &= M_2 - 2\mu(M_1) + \mu^2(M_0). \end{aligned}$$

Since $M_1 = \mu$, and $M_0 = 1$ then:

$$\begin{aligned} \text{Var} &= M_2 - 2\mu(M_1) + \mu^2(M_0) \\ &= M_2 - 2\mu(\mu) + \mu^2(1) \\ &= M_2 - 2\mu^2 + \mu^2 \\ &= M_2 - \mu^2 \\ &= M_2 - (M_1)^2. \end{aligned}$$

In other words, the variance for the uniform pdf is simply the second moment (M_2) minus the square of the first moment ($(M_1)^2$).

Thus, for a continuous uniform random variable x on the interval $[a, b]$:

$$\begin{aligned}\text{Var} &= M_2 - (M_1)^2 \\ &= \frac{a^2 + ab + b^2}{3} - \left(\frac{a+b}{2}\right)^2 = \frac{(a-b)^2}{12}.\end{aligned}$$

As it turns out, most of the usual measures by which we describe random distributions (mean, variance,...) can be expressed as functions of the moments. [We will explore this in a bit more detail later in this appendix.]

B.2. Transformations of random variables and the Delta method

OK – that’s fine. If the pdf is specified, we can use the method of moments to formally derive the mean and variance of the distribution. But, what about functions of random variables having poorly specified or unspecified distributions? Or, situations where the pdf is not easily defined?

In such cases, we may need other approaches. We introduce one such approach (the Delta method) in this Appendix, by first considering the case of a simple linear transformation of a random normal distribution.

Let

$$X_1, X_2, \dots \sim \mathcal{N}(10, \sigma^2 = 2).$$

In other words, $\{X_1, X_2, \dots\}$ are random deviates drawn from a normal distribution with a mean of $\mu = 10$, and a variance of $\sigma^2 = 2$. Now, consider some *transformations* of these random values. You might recall from some earlier statistics class that *linearly* transformed normal random variables are themselves normally distributed. Consider for example, $X_i \sim \mathcal{N}(10, 2)$ – which we then linearly transform to Y_i , such that $Y_i = 4X_i + 3$ (i.e., take each random deviate X_i , and linearly transform it multiplying it by 4, and adding 3 to it to generate the corresponding transformed variable Y_i).

Now, recall that for real scalar constants a and b (as, say, might represent the ‘slope and intercept’ in a simple linear model) we can write the expectation (E) and variance of the linear transformation as:

- i) $E(a) = a, E(aX + b) = aE(X) + b$
- ii) $\text{Var}(a) = 0, \text{Var}(aX + b) = a^2\text{Var}(X).$

[begin sidebar](#)

refresher on the algebra of expectations & variances...

Formally, the expectation operator $E[\cdot]$ satisfies the *linearity property*:

$$E[aX + b] = E[aX] + E[b].$$

Since this expectation (LHS) is a *linear operator*, then

- **scaling property:** $E[aX] = aE[X]$ because a constant multiplier (a can be factored out)
- **expectation of a constant:** $E[b] = b$ since the expected value of a constant is the constant itself

Thus,

$$\begin{aligned}E[aX + b] &= E[aX] + E[b] \\ &= aE[X] + b.\end{aligned}$$

In simple terms, (i) the expectation measures the ‘average’ value of a random variable; (ii) scaling a random variable by a scales the mean by a ; adding a constant b shifts the mean by b , but does not affect the variance.

Variance is defined (in terms of expectations) as

$$\text{Var}(Y) = E[Y^2] - (E[Y])^2,$$

where for this example, $Y = aX + b$. From above,

$$E[Y] = E[aX + b] = aE[X] + b,$$

while the expectation for Y^2 is:

$$\begin{aligned} E[Y^2] &= E[(aX + b)^2] \\ &= E[a^2X^2 + 2abX + b^2]. \end{aligned}$$

This can be re-written (using the linearity of expectations, above) as:

$$E[Y^2] = a^2E[X^2] + 2abE[X] + b^2.$$

Since

$$\begin{aligned} \text{Var}(X) &= E[X^2] - (E[X])^2 \\ E[X^2] &= \text{Var}(X) + (E[X])^2. \end{aligned}$$

Substituting this expression for $E[X^2]$ into the expression for $E[Y^2]$ (above):

$$\begin{aligned} E[Y^2] &= a^2(\text{Var}(X) + (E[X])^2) + 2abE[X] + b^2 \\ &= a^2\text{Var}(X) + a^2(E[X])^2 + 2abE[X] + b^2. \end{aligned}$$

Thus, following a few substitutions, algebraic rearrangements, and (ultimately) cancelation of some common terms:

$$\begin{aligned} \text{Var}(Y) &= E[Y^2] - (E[Y])^2 \\ &= \left(a^2\text{Var}(X) + a^2(E[X])^2 + 2abE[X] + b^2 \right) - (aE[X] + b)^2 \\ &= \left(a^2\text{Var}(X) + \cancel{a^2(E[X])^2} + \cancel{2abE[X]} + \cancel{b^2} \right) - \left(\cancel{a^2(E[X])^2} + \cancel{2abE[X]} + \cancel{b^2} \right) \\ &= a^2\text{Var}(X), \end{aligned}$$

which is the expression we started with on the previous page.

end sidebar

Thus, given $X_i \sim \mathcal{N}(10, 2)$ and the linear transformation $Y_i = 4X_i + 3$, we can write:

$$Y \sim \mathcal{N}([4(10) + 3 = 43], [(4^2)(2)]) = \mathcal{N}(43, 32).$$

Here is a simple ‘proof by demonstration’, using a few lines of **R** code:

```
X <- rnorm(10000, sd=sqrt(2))
Y <- 4*X+3
mean(Y)
42.982
var(Y)
32.121
```

An important point is that some transformations of the normal distribution are close to normal (i.e., are linear) and some are not. Since linear transformations of random normal values are also normal, it seems reasonable to conclude that *approximately* linear transformations (over some range – and this important) of random normal data should also be *approximately* normal.

Let's generalize a bit. Let $X \sim \mathcal{N}(\mu, \sigma^2)$, and let $Y = g(X)$, where g is some transformation of X (in the previous example, $g(X) = 4X + 3$ applied to $X \sim \mathcal{N}(10, 2)$). It is hopefully relatively intuitive that the closer $g(X)$ is to linear over the likely range of X (i.e., within 3 or so standard deviations of μ), the closer $Y = g(X)$ will be to 'normally distributed'. From calculus, we recall that if you look at any differentiable function over a narrow enough region, the function appears approximately linear. The approximating line is the tangent line to the curve, and its slope is the derivative of the function.

Since most of the mass (i.e., most of the random values) of X is concentrated around μ (since X is normally distributed), let's figure out the tangent line at μ , using two different methods. First, we know that the tangent line passes through $(\mu, g(\mu))$ – because when $X = \mu$, $Y = g(\mu)$ – and that its slope is $g'(\mu)$ (we use the 'prime' notation, g' , to indicate the first derivative of the function g). Thus, the equation of the tangent line is $Y = g'(\mu)X + b$ for some b . Replacing (X, Y) with the known point $(\mu, g(\mu))$, we find $g(\mu) = g'(\mu)\mu + b$ and so $b = g(\mu) - g'(\mu)\mu$. Thus, the equation of the tangent line to the transformation function, evaluated at the mean μ , is $Y = g'(\mu)X + g(\mu) - g'(\mu)\mu = g(\mu) + g'(\mu)(X - \mu)$.

Now for the big conceptual + mechanical step: we can derive an *approximation* to the same tangent line by using a *Taylor series expansion* of $g(x)$ (to first order) around $X = \mu$:

$$\begin{aligned} Y &= g(\mu) \\ &\approx g(\mu) + g'(\mu)(X - \mu) + \epsilon, \end{aligned}$$

which (after dropping the random error term ϵ) is exactly the same as what we derived above by 'logic + algebra'.

At this point you might be asking yourself 'so what?'.* Well, suppose that $X \sim \mathcal{N}(\mu, \sigma^2)$ and $Y = g(X)$, where $g'(\mu) \neq 0$. Then, whenever the tangent line (derived earlier) is approximately correct over the likely range of X (i.e., if the transformed function is approximately linear over the likely range of X), then the transformation $Y = g(X)$ will have an approximate normal distribution. Moments of that approximate normal distribution may be found using the usual rules for linear transformations of normals (which was reviewed in the -sidebar- a few pages back).

Thus, to first-order:

$$\begin{aligned} E(Y) &\approx g'(\mu)\mu + g(\mu) - g'(\mu)\mu \\ &= g(\mu), \\ \text{Var}(Y) &\approx \text{Var}(g(X)) = (g(X) - g(\mu))^2 \\ &= (g'(\mu)(X - \mu))^2 \\ &= (g'(\mu))^2 (X - \mu)^2 \\ &= (g'(\mu))^2 \text{Var}(X). \end{aligned}$$

In other words, for linear transformations, the expectation (mean), the first-order approximation is simply the transformed mean calculated for the original distribution. For the first-order approximation to the variance, we take the derivative of the transformed function with respect to the parameter, square

* Some of you might also be asking yourself 'what the heck is a Taylor series expansion?'. If so, see the next -sidebar-.

it, and multiply it by the estimated variance of the untransformed parameter.

These first-order approximations to the expectation and variance of a transformed parameter are usually referred to as the *Delta method*.*

[begin sidebar](#)

Taylor series expansions?

Briefly, the *Taylor series* is a power series expansion of an infinitely differentiable real (or complex) function defined on an open interval around some specified point.[†] For example, a one-dimensional Taylor series is an expansion of a real function $f(x)$ about a point $x = a$ over the interval $(a - r, a + r)$, is given as:

$$f(x) \approx f(a) + \frac{f'(a)(x-a)}{1!} + \frac{f''(a)(x-a)^2}{2!} + \dots,$$

where $f'(a)$ is the first derivative of f with respect to a , $f''(a)$ is the second derivative of f with respect to a , and so on.

For example, suppose the function is $f(x) = e^x$. The convenient fact about this function is that all its derivatives are equal to e^x as well (i.e., $f(x) = e^x$, $f'(x) = e^x$, $f''(x) = e^x$, \dots). In particular, $f^{(n)}(x) = e^x$ so that $f^{(n)}(0) = 1$. This means that the coefficients of the Taylor series are given by:

$$a_n = \frac{f^{(n)}(0)}{n!} = \frac{1}{n!},$$

and so the Taylor series is given by:

$$1 + x + \frac{x^2}{2} + \frac{x^3}{6} + \frac{x^4}{24} + \dots + \frac{x^n}{n!} + \dots = \sum_{n=0}^{\infty} \frac{x^n}{n!}.$$

The primary utility of such a power series in simple application is that differentiation and integration of power series can be performed term by term and is hence particularly (or, at least relatively) easy. In addition, the (truncated) series can be used to compute function values approximately.

Now, let's look at an example of the "fit" of a Taylor series to a familiar function, given a certain number of terms in the series. For our example, we'll expand the function $f(x) = e^x$, at $x = a = 0$, on the interval $[a - 2, a + 2]$, for $n = 0, n = 1, n = 2, \dots$ (where n is the number of terms in the series). For $n = 0$, the Taylor expansion is a scalar constant (1):

$$f(x) \approx 1,$$

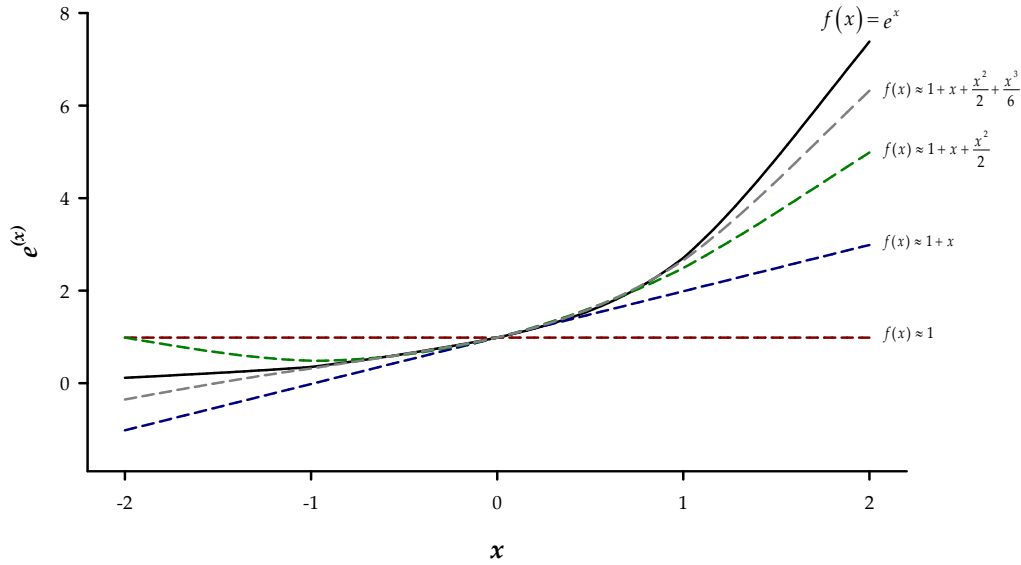
which we anticipate to be a poor approximation to the function $f(x) = e^x$ at any point. The relationship between 'the number of terms', and the 'fit' of the Taylor series expansion to the function $f(x) = e^x$, is shown clearly in the figure at the top of the next page. The solid black line in the figure is the function $f(x) = e^x$, evaluated over the interval $[-2, 2]$. The dashed lines represent different orders (i.e., number of terms) in the expansion. The red dashed line represents the 0th order expansion, $f(x) \approx 1$, the blue dashed line represents the 1st order expansion, $f(x) \approx 1 + x$, and so on.

We see in the figure that when we add more terms (i.e., use a higher-order series), the fit gets progressively better. Often, for 'nice, smooth' functions (i.e., those nearly linear at the point of interest), we don't need many terms at all. For this example, the 3rd order expansion ($n = 4$) yields a relatively good approximation to the function over much of the interval $[-2, 2]$.

Another example – suppose the function of interest is $f(x) = (x)^{1/3}$ (i.e., $f(x) = \sqrt[3]{x}$). Suppose we're interested in $f(x) = (x)^{1/3}$ where $x = 27$ (i.e., $f(27) = \sqrt[3]{27}$). Now, it is straightforward to show that

* For an interesting review of the history of the Delta method, see Ver Hoef (2012).

[†] You might have heard that the 'Taylor series' is related to something called a 'Maclaurin series'. A Maclaurin series is simply a special case of the Taylor series where the expansion is around $x = 0$ (the origin).



$f(27) = \sqrt[3]{27} = 3$. But suppose we want to know $f(25) = \sqrt[3]{25}$, using a Taylor series approximation?

We recall that to first order:

$$f(x) = f(a) + f'(a)(x - a),$$

where in this case, $a = 27$ and $x = 25$. The derivative of f with respect to x for $f(a) = (a)^{1/3}$ is:

$$f'(a) = \frac{a^{-2/3}}{3} = \frac{1}{3\sqrt[3]{a^2}}.$$

Thus, using the first-order Taylor series, we write:

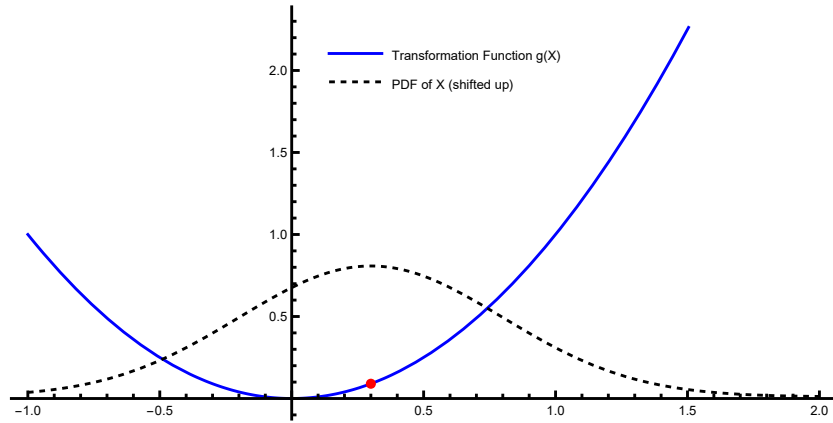
$$\begin{aligned} f(25) &\approx f(27) + f'(27)(25 - 27) \\ &= 3 + (0.037037)(-2) \\ &= 2.926. \end{aligned}$$

Clearly, 2.926 is very close to the true value of $f(25) = \sqrt[3]{25} = 2.924$. In other words, the first-order Taylor approximation works well for this particular function. As we will see later, this is not always the case, which has important implications.

end sidebar

Here is a worked example pulling all the pieces together. Suppose we have $X \sim \mathcal{N}(0.3, 0.5)$, with a transformation function $g(X) = X^2$. In other words, a simple quadratic (and therefore, non-linear) transformation of X . In the figure shown at the top of the next page, we plot the pdf for X over the interval for X of $[-1, 2]$ (black dashed line), and the transformation function $g(X)$ over the same interval (solid blue line). The red dot (•) on the transformation function corresponds to the mean of the distribution of X , $\mu = 0.3$.

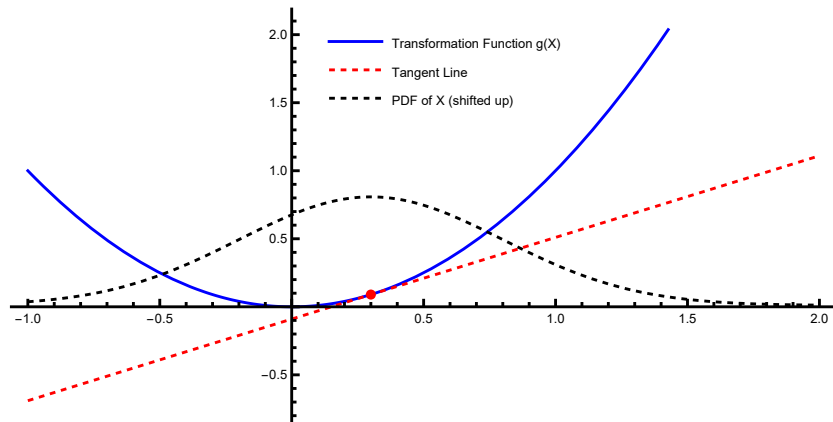
If we consider the tangent to the transformation function evaluated at the mean μ , we showed earlier that $g(\mu) = g'(\mu)\mu + b$ and so $b = g(\mu) - g'(\mu)\mu$. Since the transformation function for this example is $g(X) = X^2$, then the derivative is $g'(X) = 2X$, which evaluated at the mean is $g'(\mu) = 2(0.3) = 0.6$.



Thus, the equation of the tangent line to the function, evaluated at the mean $\mu = 0.3$, is

$$\begin{aligned} Y &= g'(\mu)X + g(\mu) - g'(\mu)\mu \\ &= 0.6X + (0.3)^2 - 0.6(0.3) \\ &= 0.6X - 0.09. \end{aligned}$$

In the following, we add this tangent line (red dashed line) evaluated at the mean μ to the plot.



Now, let's consider a first-order Taylor expansion of the transformation function $g(X)$ at $X = \mu$. Earlier, we showed that

$$\begin{aligned} Y &= g(X) \\ &\approx g'(\mu)X + g(\mu) - g'(\mu)\mu + \epsilon. \end{aligned}$$

And, recall that if we drop the error term ϵ , we end up with exactly the same equation – meaning, the Taylor approximation to first order is *exactly* the same as the tangent line to the function evaluated at μ . So, the Taylor approximation gives us the tangent line, from which we can derive the expectation and (perhaps of more interest) the variance of the transformed data.

B.3. Transformations of one variable

OK, enough background for now. Let's see some applications. Let's check the Delta method out in a few cases where we (probably) already know the answer.

Assume we have an estimate of density \hat{D} and its conditional sampling variance, $\widehat{\text{Var}}(\hat{D})$. We want to multiply this by some constant c to make it comparable with other values from the literature. Thus, we want $\hat{D}_s = g(D) = c\hat{D}$, and $\widehat{\text{Var}}(\hat{D}_s)$.

The Delta method gives:

$$\begin{aligned}\widehat{\text{Var}}(\hat{D}_s) &\approx (g'(D))^2 \hat{\sigma}_D^2 \\ &= \left(\frac{\partial \hat{D}_s}{\partial \hat{D}} \right)^2 \cdot \widehat{\text{Var}}(\hat{D}) \\ &= c^2 \cdot \widehat{\text{Var}}(\hat{D}),\end{aligned}$$

which we know to be true for the variance of a random variable multiplied by a real constant.

For example, suppose we take a large random normal sample, with mean $\mu = 1.5$, and true sample variance of $\sigma^2 = 0.25$. We are interesting in approximating the variance of this sample multiplied by some constant, $c = 2.25$. From the preceding, we expect that the variance of the transformed sample is approximated as:

$$\begin{aligned}\widehat{\text{Var}}(\hat{D}_s) &\approx c^2 \cdot \widehat{\text{Var}}(\hat{D}) \\ &= 2.25^2(0.25) \\ &= 1.265625.\end{aligned}$$

The following snippet of **R** code simulates this situation – the variance of the (rather large) random normal sample (**sample**), multiplied by the constant $c = 2.25$ (where the transformed sample is **trans_sample**) is 1.265522, which is quite close to the approximation from the Delta method (1.265625).

```
> sample <- rnorm(1000000,1.5,0.5)
> c <- 2.25
> trans_sample <- c*sample
> var(trans_sample)
[1] 1.265522
```

Another example of much the same thing – consider a known number of N harvested fish, with an average weight ($\hat{\mu}_w$) and variance. If you want an estimate of total biomass (B), then $\hat{B} = N \cdot \hat{\mu}_w$ and applying the Delta method, the variance of \hat{B} is approximated as $N^2 \cdot \widehat{\text{Var}}(\hat{\mu}_w)$.

So, if there are $N = 100$ fish in the catch, with an average mass $\hat{\mu}_w = 20$ pounds, with an estimated variance of $\hat{\sigma}_w^2 = 1.44$, then by the Delta method, the approximate variance of the total biomass $\hat{B} = (100 \times 20) = 2,000$ is:

$$\begin{aligned}\widehat{\text{Var}}(\hat{B}) &\approx N^2 \cdot \widehat{\text{Var}}(\hat{\mu}_w) \\ &= (100)^2(1.44) \\ &= 14,400.\end{aligned}$$

The following snippet of **R** code simulates this particular example – the estimated variance of the (rather large) numerically simulated sample (**biomass**), is very close to the approximation from the Delta method (14,400).

```
N <- 100 # size of sample of fish

mu <- 20; # average mass of fish
v <- 1.44; # variance of same...

reps <- 1000000 # number of replicate samples desired to generate dist of biomass

# set up replicate samples - recall that biomass is the product of N and
# the 'average' mass, which is sampled from a rnorm with mean mu and variance v
biomass <- replicate(reps,N*rnorm(1,mu,sqrt(v)));

# output result from the simulated biomass data
print(var(biomass));
[1] 14399.1
```

One final example – you have some parameter θ , which you transform by dividing it by some constant c . Thus, by the Delta method:

$$\widehat{\text{Var}}\left(\frac{\hat{\theta}}{c}\right) \approx \left(\frac{1}{c}\right)^2 \cdot \widehat{\text{Var}}(\hat{\theta}).$$

So, using a normal probability density function, for $\theta = \mu = 1.8$, and $\sigma_{\theta}^2 = 1.5$, where the constant $c = 3.14159$, then

$$\begin{aligned} \widehat{\text{Var}}\left(\frac{\hat{\theta}}{c}\right) &\approx \left(\frac{1}{3.14159}\right)^2 \cdot (1.5) \\ &= 0.15198. \end{aligned}$$

The following snippet of **R** code simulates this situation – the variance of the (rather large) random normal sample (**sample**), divided by the constant, $c = 3.14159$ (where the transformed sample is **trans_sample**) is 0.1513828, which is again quite close to the approximation from the Delta method (0.15198).

```
> sample <- rnorm(1000000,mean=1.8,sd=sqrt(1.5))
> c <- 3.14159
> sample_trans <- sample/c

> var(sample_trans)
[1] 0.1513828
```

B.3.1. A potential complication – ‘non-linear’ transformations

A final – and very important – example for transformations of single variables. The importance lies in the demonstration that the Delta method does not always work. Remember, the Delta method assumes that the transformation is approximately linear over the expected range of the parameter.

What might happen if instead the transformation is *non-linear*? Suppose one has a MLE for the mean and estimated variance for some parameter θ which is bounded random uniform on the interval $\mathcal{U}[0, 2]$. Suppose you want to transform this parameter such that:

$$\hat{g} = e^{\hat{\theta}}.$$

[This is a convenient transformation since the derivative of e^x is e^x , making the calculations very simple. Also recall from the preceding -sidebar- that the Taylor series expansion to first-order may not ‘do’ particular well with this particular transformation function.]

Now, based on the Delta method, the variance for g would be estimated as:

$$\begin{aligned}\widehat{\text{Var}}(\hat{g}) &\approx \left(\frac{\partial \hat{g}}{\partial \hat{\theta}} \right)^2 \cdot \widehat{\text{Var}}(\hat{\theta}) \\ &= (e^{\hat{\theta}})^2 \cdot \widehat{\text{Var}}(\hat{\theta}).\end{aligned}$$

Now, suppose that $\hat{\theta} = 1.0$, and $\widehat{\text{Var}}(\hat{\theta}) = 0.3\dot{3}$. Then, from the Delta method:

$$\begin{aligned}\widehat{\text{Var}}(\hat{g}) &\approx (e^{\hat{\theta}})^2 \cdot \widehat{\text{Var}}(\hat{\theta}) \\ &= (7.38906)(0.3\dot{3}) \\ &= 2.46302.\end{aligned}$$

Is this a reasonable approximation? The only way we can answer that question is if we know what the ‘true’ (correct) estimate of the variance should be.

There are two approaches we might use to come up with the ‘true’ (correct) variance: (1) analytically, or (2) by numerical simulation.

We’ll start with the formal, analytical approach, and derive the variance of g using the method of moments introduced earlier. To do this, we need to integrate the pdf (uniform, in this case) over some range. Since the variance of a uniform distribution is $(b - a)^2/12$ (as derived earlier in this appendix), and if b and a are symmetric around the mean (1.0), then we can show by algebra that given a variance of $0.3\dot{3}$, then $a = 0$ and $b = 2$ (check: $(b - a)^2/12 = (2 - 0)^2/12 = 0.3\dot{3}$).

Given a uniform distribution, the pdf is $p(\theta) = 1/(b - a)$. Thus, by the method of moments:

$$\begin{aligned}M_1 &= \int_a^b \frac{g(x)}{b - a} dx = -\frac{e^b - e^a}{a - b}, \\ M_2 &= \int_a^b \frac{g(x)^2}{b - a} dx = \left(\frac{1}{2} \right) \cdot \frac{e^{2a} - e^{2b}}{a - b}.\end{aligned}$$

Thus, $\text{Var}(E(g))$ is:

$$\begin{aligned}\text{Var}(E(g)) &= M_2 - (M_1)^2 \\ &= \left(\frac{1}{2} \right) \cdot \frac{-e^{2b} + e^{2a}}{-b + a} - \frac{(e^b - e^a)^2}{(a - b)^2}.\end{aligned}$$

If $a = 0$ and $b = 2$, then the true variance of the transformed data is given as:

$$\begin{aligned}\text{Var}(E(g)) &= M_2 - (M_1)^2 \\ &= \left(\frac{1}{2}\right) \cdot \frac{-e^{2b} + e^{2a}}{-b + a} - \frac{(e^b - e^a)^2}{(a - b)^2} \\ &= 3.19453,\end{aligned}$$

which is not particularly close to the value estimated by the Delta method (2.46302).

To confirm our analytical calculation of the variance, let’s also consider coming up with an estimate of the ‘true’ variance by numerical simulation. The steps are pretty easy: (i) simulate a large data set, (ii) transform the entire data set, and (iii) calculate the variance of the transformed data set.

For our present example, here is one way you might set this up in R:

```
> sim.data <- runif(10000000,0,2);
> transformed.data <- exp(sim.data);
> var(transformed.data);
[1] 3.19509
```

which is within Monte Carlo error of the value derived analytically, above (3.19453).

Ok, so now that we have derived the ‘true’ variance in a couple of different ways, the important question is – why the discrepancy between the ‘true’ variance of the transformed distribution (3.19453), and the first-order approximation to that variance derived using the Delta method (2.46302)?

As discussed earlier, the Delta method rests on the assumption the first-order Taylor expansion around the parameter value is effectively *linear* over the range of values likely to be encountered. Since in this example we’re using a uniform pdf, then all values between a and b are equally likely. Thus, we might anticipate that as the interval between a and b gets smaller, then the approximation to the variance (which will clearly decrease) will get better and better (since the smaller the interval, the more likely it is that the transformation function is approximately linear over that range).

For example, if $a = 0.5$ and $b = 1.5$ (same mean of 1.0), then the true variance of θ will be 0.083. Thus, by the Delta method, the estimated variance of g will be 0.61575, while by the method of moments (which is exact), the variance will be 0.65792. Clearly, the proportional difference between the two values (truth vs. estimate) has declined markedly. But, we achieved this ‘improvement’ by artificially reducing the true variance of the untransformed variable θ . Obviously, we can’t do this in general practice.

So, what are the practical options? One possible solution is to try using a higher-order Taylor series approximation. By including higher-order terms, we should achieve a better ‘fit’ to the transformation function (see the preceding -sidebar-), and thus a better approximation to the variance.

For our present example, the second-order approximation to the variance can be derived (see the -sidebar- starting on p. B-15) as:

$$\text{Var}(g(X)) \approx g'(\mu_X)^2 E[(X - \mu_X)^2] + \left(\frac{1}{4}\right) g''(\mu_X)^2 E[(X - \mu_X)^4].$$

Important technical note: Since $E[(X - \mu_X)^2] \equiv \text{Var}(X)$, you might be tempted to assume that $E[(X - \mu_X)^4]$ (the right-most term, above) is simply $(\text{Var}(X))^2$.^{*} Alas, this would be incorrect. In fact, $E[(X - \mu_X)^4]$ is what is known as the ‘fourth *central* moment of the pdf’.

^{*} Since, for example, if $x = (a - b)^2$, where a and b are real scalars, then $(a - b)^4 = [(a - b)^2]^2 = x^2$.

What is a *central* moment? There are two different types of moments for a continuous pdf: *raw*, and *central*. A *raw* moment describes the overall shape of a distribution, while *central* moments measure the expected values of powers of deviations from the mean. Put another way, *raw* moments are absolute measures of the ‘shape’ of a distribution, while *central* moments adjust for the mean and focus on spread and shape relative to it. At the start of this appendix, we derived the moments M_0, M_1, \dots for the uniform pdf $\mathcal{U}[a, b]$. At the time, we didn’t specify that these were in fact ‘raw’ moments. We showed that we could use these raw moments to derive the variance of the distribution, as $\text{Var} = M_2 - (M_1)^2$. In fact, all of the familiar central moments (things like variance, skewness, kurtosis,...) are derived from the raw moments (for example, see derivation of $\mu_3 = E[(X - \mu_X)^3]$ in Appendix B.2).

Since central moments focus on deviations relative to the mean, we can generally express the central moment of order n as $\mu_n = E[(X - \mu_X)^n]$:

$$\begin{aligned}\mu_n &= E[(X - \mu_X)^n] \\ &= \sum_{k=0}^n \binom{n}{k} (-\mu_X)^{n-k} M_k,\end{aligned}$$

where μ_n is the n th central moment, M_k is the k th raw moment, μ_X is the mean of the distribution of X , and $\binom{n}{k}$ is the binomial coefficient given n and k .

So, the fourth central moment, $\mu_4 = E[(X - \mu_X)^4]$, is given for $\mathcal{U}[a, b]$ as

$$\begin{aligned}\mu_4 &= \sum_{k=0}^4 \binom{4}{k} (-\mu_X)^{4-k} M_k \\ &= \binom{4}{0} (-\mu_X)^4 M_0 + \binom{4}{1} (-\mu_X)^3 M_1 + \binom{4}{2} (-\mu_X)^2 M_2 + \binom{4}{3} (-\mu_X)^1 M_3 + \binom{4}{4} (-\mu_X)^0 M_4 \\ &= (1)(-\mu_X)^4 M_0 + (4)(-\mu_X)^3 M_1 + (6)(-\mu_X)^2 M_2 + (4)(-\mu_X) M_3 + (1)(M_4).\end{aligned}$$

Since $M_0 = 1$, and $M_1 = \mu_X$, then after some substitutions and a bit of algebra, the fourth central moment for $\mathcal{U}[a, b]$ is

$$\mu_4 = E[(X - \mu_X)^4] = M_4 - 4M_1M_3 + 6M_1^2M_2 - 3M_1^4.$$

For our example, based on $\mathcal{U}[0, 2]$, $M_0 = 1$, $M_1 = 1$, with (skipping the calculations) $M_2 = 1.\dot{3}$, $M_3 = 2$ and $M_4 = 3.2$. Substituting these into our equation for μ_4 (above):

$$\begin{aligned}\mu_4 &= E[(X - \mu_X)^4] = M_4 - 4M_1M_3 + 6M_1^2M_2 - 3M_1^4 \\ &= 3.2 - 4(1)(2) + 6(1)^2(1.\dot{3}) - 3(1)^4 \\ &= 0.2.\end{aligned}$$

Since the derivative of e^x is the same regardless of the order of the derivative (i.e., $g' = g'' = \dots = e^x$), and since $\mu_X = 1.0$, $\text{Var}(X) = 0.3\dot{3}$ for $X \sim \mathcal{U}[0, 2]$, then to second-order:

$$\begin{aligned}\text{Var}(g(X)) &\approx g'(\mu_X)^2 E[(X - \mu_X)^2] + \frac{1}{4} g''(\mu_X)^2 E[(X - \mu_X)^4] \\ &= e^2 \cdot (0.3\dot{3}) + \frac{1}{4} e^2 \cdot (0.2) = 2.83247,\end{aligned}$$

which is ~50% closer to the true variance (3.193714) than the first-order approximation (2.46302), but is still not a particularly good estimate.

Key take-home message \Rightarrow if your transformation is strongly non-linear around the mass of the data, the Delta method might not yield robust estimates of the variance.

begin sidebar

deriving higher order approximating equations...

In the preceding, we considered a second-order approximation to the variance of the transformed data. Where did that particular approximating equation come from? How would we derive expression for even higher-order expansions?

Recall from earlier in this Appendix that the *first-order* approximation follows from a Taylor series expansion of $g(X)$ around $E[X]$:

$$g(X) \approx g(\mu_X) + g'(\mu_X)(X - \mu_X).$$

Taking variance on both sides:

$$\text{Var}(g(X)) \approx \text{Var}(g'(\mu_X)(X - \mu_X)).$$

Since $g'(\mu_X)$ is a constant, we apply the property $\text{Var}(aX) = a^2\text{Var}(X)$:

$$\text{Var}(g(X)) \approx g'(\mu_X)^2 \text{Var}(X).$$

Expressing variance in terms of expectation,

$$\text{Var}(g(X)) \approx (g'(E[X]))^2 E[(X - E[X])^2].$$

The *second-order* Delta method approximation follows the same process – i.e., using the Taylor expansion of $g(X)$ around $E[X]$ – but including a second-order term:

$$g(X) \approx g(\mu_X) + g'(\mu_X)(X - \mu_X) + \frac{1}{2}g''(\mu_X)(X - \mu_X)^2.$$

Taking the expectation:

$$E[g(X)] \approx g(\mu_X) + g'(\mu_X)E[X - \mu_X] + \frac{1}{2}g''(\mu_X)E[(X - \mu_X)^2].$$

Since $E[X - \mu_X] = 0$, this simplifies to:

$$E[g(X)] \approx g(\mu_X) + \frac{1}{2}g''(\mu_X)E[(X - \mu_X)^2].$$

The variance of $g(X)$ is given by:

$$\text{Var}(g(X)) = E[(g(X) - E[g(X)])^2].$$

Using the second-order expansion (above):

$$g(X) - E[g(X)] \approx g'(\mu_X)(X - \mu_X) + \frac{1}{2}g''(\mu_X)[(X - \mu_X)^2 - E[(X - \mu_X)^2]].$$

Squaring both sides and expanding the square on the RHS:

$$\begin{aligned} (g(X) - E[g(X)])^2 &\approx \left[g'(\mu_X)(X - \mu_X) + \frac{1}{2}g''(\mu_X)[(X - \mu_X)^2 - E[(X - \mu_X)^2]] \right]^2 \\ &= g'(\mu_X)^2(X - \mu_X)^2 \\ &\quad + g'(\mu_X)g''(\mu_X)(X - \mu_X)[(X - \mu_X)^2 - E[(X - \mu_X)^2]] \end{aligned}$$

$$+ \frac{1}{4} g''(\mu_X)^2 [(X - \mu_X)^2 - E[(X - \mu_X)^2]]^2.$$

Taking expectations:

$$\begin{aligned} \text{Var}(g(X)) &\approx g'(\mu_X)^2 E[(X - \mu_X)^2] \\ &\quad + g'(\mu_X) g''(\mu_X) E[(X - \mu_X)^3] \\ &\quad + \frac{1}{4} g''(\mu_X)^2 E[(X - \mu_X)^4]. \end{aligned}$$

A key (and useful) observation is that for *symmetric* distributions like the normal or uniform, the odd central moments 3 and higher (3, 5, 7,...) are zero (see Addendum B.2), and can be dropped from the expression. So, for this example:

$$E[(X - \mu_X)^3] = 0.$$

Thus, removing the third-moment term we end up with the second-order equation we applied in our example:

$$\text{Var}(g(X)) \approx g'(\mu_X)^2 E[(X - \mu_X)^2] + \frac{1}{4} g''(\mu_X)^2 E[(X - \mu_X)^4].$$

General k -th Order Approximation

If we were so inclined*, we can generalize the preceding as follows. Using a Taylor expansion:

$$g(X) = g(\mu_X) + g'(\mu_X)(X - \mu_X) + \frac{1}{2} g''(\mu_X)(X - \mu_X)^2 + \cdots + \frac{1}{k!} g^{(k)}(\mu_X)(X - \mu_X)^k.$$

Taking variance and simplifying, we obtain:

$$\text{Var}(g(X)) \approx \sum_{k=1}^{\infty} \frac{g^{(k)}(\mu_X)^2}{(k!)^2} E[(X - \mu_X)^{2k}].$$

As noted, for a symmetric distribution, all odd central moments vanish (Addendum B.2), simplifying the expression to:

$$\text{Var}(g(X)) \approx \sum_{k=1}^{\infty} \frac{g^{(2k)}(\mu_X)^2}{(2k)!^2} E[(X - \mu_X)^{4k}].$$

For symmetric distributions, only even derivatives of $g(X)$ contribute, and the first correction beyond the standard first-order Delta method occurs at order 2 (involving the fourth central moment).

For small k , this yields:

- **First-order approximation:**

$$\text{Var}(g(X)) \approx g'(\mu_X)^2 E[(X - \mu_X)^2].$$

- **Second-order approximation:**

$$\text{Var}(g(X)) \approx g'(\mu_X)^2 E[(X - \mu_X)^2] + \frac{1}{4} \cdot g''(\mu_X)^2 E[(X - \mu_X)^4].$$

- **Third-order correction:**

$$\text{Var}(g(X)) \approx g'(\mu_X)^2 E[(X - \mu_X)^2] + \frac{g''(\mu_X)^2}{4} E[(X - \mu_X)^4] + \frac{g'''(\mu_X)^2}{36} E[(X - \mu_X)^6].$$

end sidebar

* And perhaps slightly masochistic...

B.4. Delta method & the expectation of transformed data

Consider the following situation. Suppose you are interested in simulating some data on the logit scale, where variation around the mean is normal (so, you're going to simulate logit-normal data). Suppose the mean of some parameter on the real probability scale is $\theta = 0.3$. Transformed to the logit scale, the mean of the sample you're going to simulate would be $\log(\theta/(1-\theta)) = -0.8472979$. So, you want to simulate some normal data, with some specified variance, on the logit scale, centered on $\mu_{\text{logit}} = -0.8472979$.

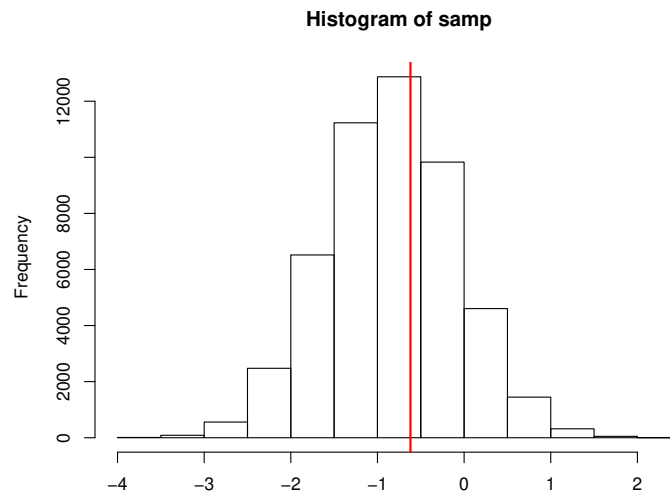
Here, using **R**, we generate a vector (which we've called 'samp', below) of 100,000 logit-normal deviates, with a $\mu_{\text{logit}} = -0.8472979$, and a standard deviation of $\sigma_{\text{logit}} = 0.75$ (corresponding to a variance of $\sigma_{\text{logit}}^2 = 0.5625$).

```
> samp <- rnorm(100000, -0.8472979, 0.75);
```

If we check the mean and variance of our random sample, we find they're quite close to the true parameters used in simulating the data (perhaps not surprising given the size of the simulated sample).

```
> mean(samp)
[1] -0.8451009
> var(samp)
[1] 0.5618909
```

If we plot a histogram of the simulated data (below), we see a symmetrical distribution centered around the true mean $\mu_{\text{logit}} = -0.8472979$ (vertical red line).



What is the variance of the back-transformed estimate of the mean, on the real probability scale? We know from what we've covered so far that if we try to calculate the variance of the back-transform of these data from the logit scale \rightarrow real probability scale, by simply taking the back transform of the estimated variance $\hat{\sigma}_{\text{logit}}^2 = 0.5618909$, we'll get the incorrect answer. If we do that, we would get

$$\frac{e^{0.5618909}}{1 + e^{0.5618909}} = 0.6368899.$$

How can we confirm our developing intuition that this value is incorrect? Well, if we simply back-transform the entire random sample, and then calculate the variance of this back transformed sample (which we call **'back'**) directly,

```
> expit=function(x) exp(x)/(1+exp(x));
> back <- expit(samp)
> var(back)
[1] 0.02211352
```

we get a value which, as we might have expected, isn't remotely close to the value of 0.6368899 we obtained by simply back-transforming the variance estimate.

Of course, we know by now we should have used the Delta method here. First, we recall that the back-transform f from the logit \rightarrow to the real scale is:

$$f = \frac{e^{\theta}}{1 + e^{\theta}}.$$

Then, we apply the Delta method as:

$$\begin{aligned}\widehat{\text{Var}}(f) &\approx \left(\frac{\partial f}{\partial \hat{\theta}} \right)^2 \times \widehat{\text{Var}}(\hat{\theta}) \\ &= \left(\frac{e^{\hat{\theta}}}{(1 + e^{\hat{\theta}})^2} \right)^2 \times \widehat{\text{Var}}(\hat{\theta}) \\ &= \left(\frac{e^{-0.8451009}}{(1 + e^{-0.8451009})^2} \right)^2 \times 0.5618909 = 0.02482293,\end{aligned}$$

which is quite close to the value we derived (above) by calculating the variance of the entire back-transformed sample directly (0.02211352).

However, the main point we want to cover here is applying the Delta method to other moments. Specifically, the mean. Recall that the mean from our logit-normal sample was -0.8451009 . Can we simply back-transform this mean from the logit \rightarrow real probability scale? In other words,

$$\frac{e^{-0.8451009}}{1 + e^{-0.8451009}} = 0.3004616.$$

Now, compare this value to the mean of the entire back-transformed sample:

```
> mean(back)
[1] 0.3199998
```

You might think that the two values (0.3004616, 0.3199998) are 'close enough for government work' (although the difference is roughly 6%), but since we don't work for the government, let's apply the Delta method to generate a correct approximation to the back-transformed mean.

First, recall that the transformation function f (from logit \rightarrow real) is

$$f = \frac{e^{\theta}}{1 + e^{\theta}}.$$

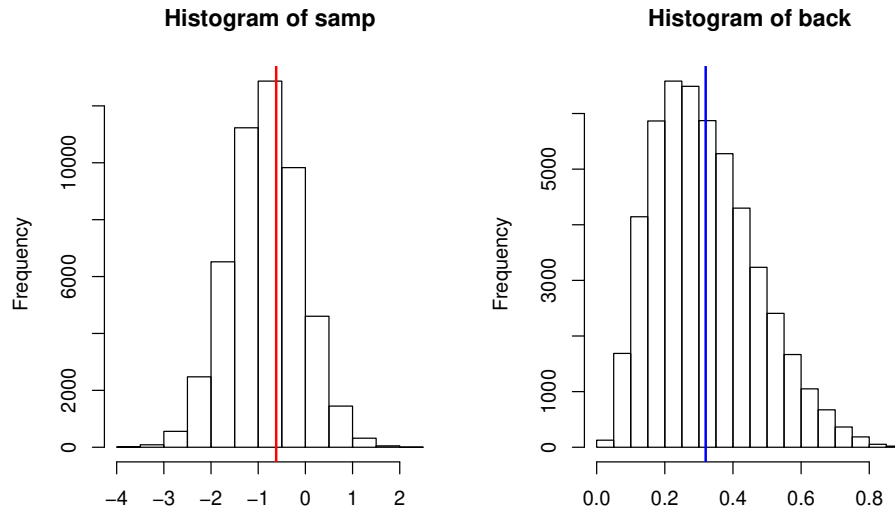
Next, remember that the Delta method as we've been applying generally (and in the preceding for the variance) it is based on the first-order Taylor series approximation.

What is the first-order Taylor series expansion for f , if $\theta = \mu$? In fact, it is simply:

$$\frac{e^\mu}{1 + e^\mu} + O((\theta - \mu)^2),$$

where the error term $O((\theta - \mu)^2)$ is the asymptotic bound of the growth of the error. But, more to the point, the first-order approximation is basically our back-transformation, with some (possibly a lot) of error added. In fact, we might expect this error term to be increasingly important if the assumptions under which the first-order approximation applies are strongly violated. In particular, if the transformation function is highly non-linear over the range of values being examined.

Do we have such a situation in the present example? Compare the histograms of our simulated data, on the logit (**samp**) and back-transformed real scales (**back**), respectively, shown below:



Note that the mean of the back-transformed distribution (vertical blue line) is somewhat to the right of the mass of the distribution, which is fairly asymmetrical. This suggests that the back-transformation might be sufficiently non-linear that we need to use a higher-order Taylor series approximation.

If you do the math (which isn't that difficult), the second-order approximation is given as

$$\frac{e^\mu}{(1 + e^\mu)} + \frac{e^\mu}{(1 + e^\mu)^2}(\theta - \mu) + O((\theta - \mu)^2).$$

Now, while this might look 'complicated', the key is to remember that we're dealing with 'expectations'. What is the expectation of $(\theta - \mu)$? In this situation, θ is a random variable – where each estimated mean from a set of replicated data sets on the logit scale represents θ , and μ is the overall parametric mean. We know from even the most basic statistics class that the expectation of the difference of a random variable X_i from the mean of the set of random variables, \bar{X} , is 0 (i.e., $E(X_i - \bar{X}) = 0$).

By the same logic, then, the expectation of $E(\theta - \mu) = 0$. And, anything multiplied by 0 is 0, so, after dropping the error term, our second-order approximation reduces to

$$\frac{e^\mu}{(1 + e^\mu)} + \cancel{\frac{e^\mu}{(1 + e^\mu)^2}(\theta - \mu)} = \frac{e^\mu}{(1 + e^\mu)},$$

which brings us right back to our standard first-order approximation.

What about a third-order approximation? After a bit more math, we end up with

$$\frac{e^\mu}{(1 + e^\mu)} + \frac{e^\mu}{(1 + e^\mu)^2}(\theta - \mu) - \frac{1}{2} \frac{e^\mu(e^\mu - 1)}{(1 + e^\mu)^3}(\theta - \mu)^2 + O((\theta - \mu)^3).$$

Again, the expectation for $E(\theta - \mu) = 0$. So, that term drops out:

$$\frac{e^\mu}{(1 + e^\mu)} + \cancel{\frac{e^\mu}{(1 + e^\mu)^2}(\theta - \mu)} - \frac{1}{2} \frac{e^\mu(e^\mu - 1)}{(1 + e^\mu)^3}(\theta - \mu)^2$$

What about for the term $E(\theta - \mu)^2$? Look closely – ‘*variate minus mean, squared*’. Look familiar? It should – it’s the variance! So, $E(\theta - \mu)^2 = \hat{\sigma}^2$.

Thus, after dropping the error term, our third-order Delta approximation to the mean is given as

$$\frac{e^\mu}{(1 + e^\mu)} - \frac{1}{2} \frac{e^\mu(e^\mu - 1)}{(1 + e^\mu)^3} \sigma^2.$$

So, given our estimate of $\hat{\mu} = -0.8451009$ and $\hat{\sigma}^2 = 0.5618909$ on the logit-scale, our third-order Delta method approximation for the expectation (mean) on the back-transformed real probability scale, using this third-order approximation is

$$\frac{e^{-0.8451009}}{(1 + e^{-0.8451009})} - \frac{1}{2} \frac{e^{-0.8451009}(e^{-0.8451009} - 1)}{(1 + e^{-0.8451009})^3} (0.5618909) = 0.3240272,$$

which is quite a bit closer to the empirical estimate of the mean derived from the entire back-transformed sample (0.3199998) than was our first attempt using the first-order approximation (0.3004616).

So, we see that the classical Delta method, which is based on a first-order Taylor series expansion of the transformed function, may not do particularly well if the function is highly non-linear over the range of values being examined. Of course, it would be fair to note that the preceding example made the assumption that the distribution was random uniform over the interval.

For most of our work with **MARK**, the interval is likely to have a near-symmetric mass around the estimate, typically β . As such, most of data, and thus the transformed data, will actually fall closer to the parameter value in question (the mean in this example) than we’ve demonstrated here. So much so, that the discrepancy between the first order ‘Delta’ approximation to the variance and the true value of the variance will likely be significantly smaller than shown here, even for a strongly non-linear transformation. We leave it to you as an exercise to prove this for yourself.

But, this point notwithstanding, it is important to be aware of the assumptions underlying the Delta method. If your transformation is non-linear, and there is considerable variation in your data, the first-order approximation may not be particularly good.

B.5. Transformations of two or more variables

We are very often interested in transformations involving more than one variable. Fortunately, there are multivariate generalizations of the Delta method. For example, the product of $(\hat{\phi}_1 \times \hat{\phi}_2 \times \hat{\phi}_3)$ we motivated this appendix with is a multivariate transformation of the 3 survival estimates.

Suppose you've estimated p different random variables X_1, X_2, \dots, X_p . In matrix notation, these variables would constitute a $(p \times 1)$ random vector \mathbf{X} :

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}, \text{ which has a corresponding mean vector } \boldsymbol{\mu} = \begin{bmatrix} EX_1 \\ EX_2 \\ \vdots \\ EX_p \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix},$$

and $(p \times p)$ variance-covariance matrix:

$$\text{Cov}(\mathbf{X}) = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Var}(X_p) \end{bmatrix}.$$

If the variables X_1, X_2, \dots, X_p are *independent*, then the off-diagonal elements (i.e., covariance terms) are all zero.

Thus, for a $(k \times p)$ matrix of constants $\mathbf{A} = a_{ij}$, the expectation of a random vector $\mathbf{Y} = \mathbf{AX}$ is given as:

$$\begin{bmatrix} EY_1 \\ EY_2 \\ \vdots \\ EY_p \end{bmatrix} = \mathbf{A}\boldsymbol{\mu},$$

with a variance-covariance matrix

$$\text{Cov}(\mathbf{Y}) = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T.$$

Using the same logic we applied in developing the Delta method for a single variable, for each x_i near μ_i , we can write:

$$\mathbf{y} = \begin{bmatrix} g_1(x) \\ g_2(x) \\ \vdots \\ g_p(x) \end{bmatrix} \approx \begin{bmatrix} g_1(\mu) \\ g_2(\mu) \\ \vdots \\ g_p(\mu) \end{bmatrix} + \mathbf{D}(x - \mu),$$

where \mathbf{D} is the matrix of partial derivatives of g_i with respect to x_j , evaluated at $(x - \mu)$.

As with the single-variable Delta method, if the variances of the X_i are small (so that with high

probability Y is near μ , such that the linear approximation is likely to be valid), then to first-order we can write:

$$\begin{bmatrix} EY_1 \\ EY_2 \\ \vdots \\ EY_p \end{bmatrix} = \begin{bmatrix} g_1(\mu) \\ g_2(\mu) \\ \vdots \\ g_p(\mu) \end{bmatrix}$$

$$\widehat{\text{Var}}(\hat{Y}) \approx \mathbf{D}\Sigma\mathbf{D}^\top.$$

In other words, to approximate the variance of some multi-variable function \mathbf{Y} , we (i) take the vector of partial derivatives of the function with respect to each parameter in turn (generally known as the *Jacobian*), \mathbf{D} , (ii) right-multiply this vector by the variance-covariance matrix, Σ , and (iii) right-multiply the resulting product by the transpose of the original vector of partial derivatives, \mathbf{D}^\top .

Note: interpretation of the variance estimated using the Delta method is dependent on the source of the variance-covariance matrix, Σ , used in the calculations. If Σ is constructed using standard ML estimates of the variances and covariances, then the resulting Delta method estimate for variance is an estimate of the ‘total’ variance, which is the sum of ‘sampling’ + ‘biological process’ variance. In contrast, if Σ is based on estimated ‘process’ variances and covariances only, then the Delta method estimate for variance is an estimate of the ‘process’ variance. Decomposition of total variance into sampling and process components is covered in detail in Appendix D.

begin sidebar

alternative algebras for Delta method

There are alternative formulations of this expression which may be more convenient to implement in some instances. When the variables $\theta_1, \theta_2, \dots, \theta_k$ (in the function, Y) are independent, then

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \text{Var}(f(\theta_1, \theta_2, \dots, \theta_k)) \\ &= \sum_{i=1}^k \text{Var}(\theta_i) \left(\frac{\partial f}{\partial \theta_i} \right)^2, \end{aligned}$$

where $\partial f / \partial \theta_i$ is the partial derivative of Y with respect to θ_i .

When the variables $\theta_1, \theta_2, \dots, \theta_k$ (in the function, Y) are **not** independent, then the covariance among the variables must be accounted for:

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \text{Var}(f(\theta_1, \theta_2, \dots, \theta_k)) \\ &= \sum_{i=1}^k \text{Var}(\theta_i) \left(\frac{\partial f}{\partial \theta_i} \right)^2 + 2 \sum_{i < j}^k \text{Cov}(\theta_i, \theta_j) \left(\frac{\partial f}{\partial \theta_i} \right) \left(\frac{\partial f}{\partial \theta_j} \right). \end{aligned}$$

end sidebar

Example (1) – variance of a product of survival probabilities

Let’s consider the application of the Delta method in estimating sampling variances of a fairly common function – the product of several parameter estimates.

From the preceding, we see that:

$$\begin{aligned}\widehat{\text{Var}}(\hat{Y}) &\approx \mathbf{D}\mathbf{\Sigma}\mathbf{D}^T \\ &= \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right] \cdot \hat{\mathbf{\Sigma}} \cdot \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right]^T,\end{aligned}$$

where Y is some linear or nonlinear function of the k parameter estimates $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$.

The first term, \mathbf{D} , on the RHS of the variance expression is a row vector containing partial derivatives of Y with respect to each of these k parameters ($\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$). The right-most term of the RHS of the variance expression, \mathbf{D}^T , is simply a transpose of this row vector (i.e., a column vector). The middle-term, $\mathbf{\Sigma}$ is simply the estimated variance-covariance matrix for the parameters.

To demonstrate the steps in the calculation, we'll use estimates from model $\{\varphi_t p.\}$ fit to the male European dipper data set. Suppose we're interested in the probability of surviving from the start of the first interval to the end of the third interval.

The estimate of this probability is easy enough:

$$\begin{aligned}\hat{Y} &= (\hat{\varphi}_1 \times \hat{\varphi}_2 \times \hat{\varphi}_3) \\ &= (0.6109350 \times 0.458263 \times 0.4960239) \\ &= 0.138871.\end{aligned}$$

So, the estimated probability of a male Dipper surviving over the first three intervals is $\sim 14\%$ (again, assuming that our time-dependent survival model is a valid model).

To derive the estimate of the variance of the product, we will also need the variance-covariance matrix for the survival estimates. You can access the matrix fairly easily in **MARK** by selecting '**Output | Specific Model Output | Variance-Covariance Matrices | Real Estimates**'.

The variance-covariance matrix for the male Dipper data, generated from model $\{\varphi_t p.\}$, as output to the default editor (e.g., Windows Notepad), is shown below:

Variance-Covariance matrix of estimates on diagonal and below, Correlation matrix of estimates above diagonal.						
	1 7	2	3	4	5	6
1	0.02243 -0.09253	-0.02638	0.00513	0.00735	0.00516	0.02379
2	-0.00039 -0.06865	0.00997	-0.02779	0.00545	0.00383	0.01765
3	0.00007 -0.05549	-0.00024	0.00724	-0.03332	0.00309	0.01427
4	0.00009 -0.07941	0.00004	-0.00023	0.00661	-0.04175	0.02042
5	0.00006 -0.05572	0.00003	0.00002	-0.00026	0.00581	-0.02857
6	0.00028 -0.25711	0.00014	0.00010	0.00013	-0.00017	0.00634
7	-0.00053 0.00146	-0.00026	-0.00018	-0.00025	-0.00016	-0.00078

The variance-covariance values are shown *below* the diagonal, whereas the standardized correlation values are *above* the diagonal. The variances are given *along* the diagonal. [Remember that for this example, we are interested in the variances and covariances among $\hat{\phi}_1 \rightarrow \hat{\phi}_3$ only, corresponding to the upper (3×3) sub-matrix of the full V-C matrix.]

However, it is **very important** to note that the V-C matrix **MARK** outputs to the editor is *rounded* to 5 significant digits. For our calculations, we need to use the full precision values.* To get those, you need to either (i) output the V-C matrix into a dBase file (which you could then open with dBase, or Excel), or (ii) copy the V-C matrix into the Windows clipboard, and then paste it into some other application. Failure to use the full precision V-C matrix will almost always lead to significant ‘rounding errors’.

The ‘full precision’ V-C matrix ($\hat{\Sigma}$) for the 3 Dipper survival estimates is shown below:

$$\hat{\Sigma} = \begin{bmatrix} \widehat{\text{Var}}(\hat{\phi}_1) & \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_2) & \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_3) \\ \widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_1) & \widehat{\text{Var}}(\hat{\phi}_2) & \widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_3) \\ \widehat{\text{Cov}}(\hat{\phi}_3, \hat{\phi}_1) & \widehat{\text{Cov}}(\hat{\phi}_3, \hat{\phi}_2) & \widehat{\text{Var}}(\hat{\phi}_3) \end{bmatrix} = \begin{bmatrix} 0.0224330125 & -0.0003945405 & 0.0000654469 \\ -0.0003945405 & 0.0099722201 & -0.0002361998 \\ 0.0000654469 & -0.0002361998 & 0.0072418858 \end{bmatrix}.$$

For this example, the transformation we’re applying to our 3 survival estimates (which we’ll call Y) is the product of the estimates (i.e., $\hat{Y} = \hat{\phi}_1 \hat{\phi}_2 \hat{\phi}_3$).

Thus, our variance estimate is given as

$$\widehat{\text{Var}}(\hat{Y}) \approx \begin{bmatrix} \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_1} \right) & \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_2} \right) & \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_3} \right) \end{bmatrix} \cdot \hat{\Sigma} \cdot \begin{bmatrix} \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_1} \right) \\ \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_2} \right) \\ \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_3} \right) \end{bmatrix}.$$

Each of the partial derivatives for \hat{Y} is easy enough to derive for this example. Since $\hat{Y} = \hat{\phi}_1 \hat{\phi}_2 \hat{\phi}_3$, then $\partial \hat{Y} / \partial \hat{\phi}_1 = \hat{\phi}_2 \hat{\phi}_3$. And so on, resulting in:

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \begin{bmatrix} \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_1} \right) & \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_2} \right) & \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_3} \right) \end{bmatrix} \cdot \hat{\Sigma} \cdot \begin{bmatrix} \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_1} \right) \\ \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_2} \right) \\ \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_3} \right) \end{bmatrix} \\ &= \begin{bmatrix} (\hat{\phi}_2 \hat{\phi}_3) & (\hat{\phi}_1 \hat{\phi}_3) & (\hat{\phi}_1 \hat{\phi}_2) \end{bmatrix} \cdot \begin{bmatrix} \widehat{\text{Var}}(\hat{\phi}_1) & \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_2) & \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_3) \\ \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_1) & \widehat{\text{Var}}(\hat{\phi}_2) & \widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_3) \\ \widehat{\text{Cov}}(\hat{\phi}_3, \hat{\phi}_1) & \widehat{\text{Cov}}(\hat{\phi}_3, \hat{\phi}_2) & \widehat{\text{Var}}(\hat{\phi}_3) \end{bmatrix} \cdot \begin{bmatrix} (\hat{\phi}_2 \hat{\phi}_3) \\ (\hat{\phi}_1 \hat{\phi}_3) \\ (\hat{\phi}_1 \hat{\phi}_2) \end{bmatrix}. \end{aligned}$$

Clearly, the approximation is getting more and more ‘impressive looking’ as we progress.

* The variance-covariance estimates **MARK** generates will occasionally depend on which optimization method you use (i.e., default, or simulated annealing), and on the starting values used to initialize the optimization. The differences are often very small (i.e., apparent only several decimal places out from zero), but you should be aware of them. For the examples presented in this Appendix, we have used the default optimization routines, and default starting values.

The resulting expression (written in piecewise fashion to make it easier to see the basic pattern) is shown below:

$$\begin{aligned}\widehat{\text{Var}}(\hat{Y}) &\approx \hat{\phi}_2^2 \hat{\phi}_3^2 [\widehat{\text{Var}}(\hat{\phi}_1)] \\ &\quad + 2\hat{\phi}_2 \hat{\phi}_3^2 \hat{\phi}_1 [\widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_2)] \\ &\quad + 2\hat{\phi}_2^2 \hat{\phi}_3 \hat{\phi}_1 [\widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_3)] \\ &\quad + \hat{\phi}_1^2 \hat{\phi}_3^2 [\widehat{\text{Var}}(\hat{\phi}_2)] \\ &\quad + 2\hat{\phi}_1^2 \hat{\phi}_3 \hat{\phi}_2 [\widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_3)] \\ &\quad + \hat{\phi}_1^2 \hat{\phi}_2^2 [\widehat{\text{Var}}(\hat{\phi}_3)].\end{aligned}$$

After substituting in our estimates for ϕ_i and the variances and covariances, our estimate for the variance of the product $\hat{Y} = (\hat{\phi}_1 \hat{\phi}_2 \hat{\phi}_3)$ is (approximately) $\widehat{\text{Var}}(Y) = 0.0025565$.

Example (2) – variance of estimate of reporting rate

In some cases animals are tagged or banded to enable estimation of a “reporting rate” – the proportion of tagged animals reported (say, to a conservation management agency), given that they were killed and retrieved by a hunter or angler (see chapter 8 for more details). Thus, N_c animals are tagged with normal (*control*) tags and, of these, R_c are recovered the first year following release. The *recovery rate* of these control animals is merely R_c/N_c and we denote this as f_c .

Another group of animals, of sample size N_r , are tagged with special *reward* tags; these tags indicate that some amount of money (say, \$50) will be given to people reporting these special tags. It is assumed that \$50 is sufficient to ensure that all such tags will be reported, thus these serve as a basis for comparison and the estimation of a reporting rate. The recovery probability for the reward tagged animals is merely R_r/N_r , where R_r is the number of recoveries of reward-tagged animals the first year following release. We denote this recovery probability as f_r .

The estimator of the *reporting rate* is a ratio of the *recovery rates* and we denote this as λ :

$$\hat{\lambda} = \frac{\hat{f}_c}{\hat{f}_r}.$$

Note that both recovery probabilities are binomials (i.e., you’re either recovered, or not). As such, we can construct variance estimators for \hat{f}_c and \hat{f}_r in the usual way:

$$\widehat{\text{Var}}(\hat{f}_c) = \frac{\hat{f}_c(1 - \hat{f}_c)}{N_c} \quad \text{and} \quad \widehat{\text{Var}}(\hat{f}_r) = \frac{\hat{f}_r(1 - \hat{f}_r)}{N_r}.$$

In this case, the samples are independent (since we assume – reasonably – that whether the tag or band is ‘control’ or ‘reward’ doesn’t actually influence the recovery probability), thus $\text{Cov}(f_c, f_r)$ and the sampling variance-covariance matrix is diagonal:

$$\begin{bmatrix} \widehat{\text{Var}}(\hat{f}_c) & 0 \\ 0 & \widehat{\text{Var}}(\hat{f}_r) \end{bmatrix}.$$

Next, we need the derivatives of λ with respect to f_c and f_r :

$$\frac{\partial \hat{\lambda}}{\partial \hat{f}_c} = \frac{1}{\hat{f}_r}, \quad \text{and} \quad \frac{\partial \hat{\lambda}}{\partial \hat{f}_r} = -\frac{\hat{f}_c}{\hat{f}_r^2}.$$

Thus,

$$\widehat{\text{Var}}(\hat{\lambda}) \approx \begin{bmatrix} \frac{1}{\hat{f}_r} & -\frac{\hat{f}_c}{\hat{f}_r^2} \end{bmatrix} \begin{bmatrix} \widehat{\text{Var}}(\hat{f}_c) & 0 \\ 0 & \widehat{\text{Var}}(\hat{f}_r) \end{bmatrix} \begin{bmatrix} \frac{1}{\hat{f}_r} \\ \frac{\hat{f}_c}{\hat{f}_r^2} \end{bmatrix}.$$

Example (3) – variance of back-transformed estimates - simple

In Chapter 6, we demonstrated how we can ‘back-transform’ from the estimate of β on the logit scale to an estimate of some parameter θ (e.g., φ or p) on the real probability scale (which is bounded $[0, 1]$). But, we’re clearly also interested in an estimate of the variance (precision) of our estimate, on both scales. Your first thought might be to simply back-transform from the link function (in our example, the logit link), to the probability scale, just as we did above. But, as discussed in chapter 6, this does not work.

For example, consider the male Dipper data. Using the logit link, we fit the time-invariant model $\{\varphi, p.\}$ to the data. Let’s consider only the estimate for $\hat{\varphi}$. The estimate for $\hat{\beta}$ for φ is 0.2648275. Thus, our estimate of $\hat{\varphi}$ on the probability scale (which is what **MARK** reports) is:

$$\hat{\varphi} = \frac{e^{0.2648275}}{1 + e^{0.2648275}} = \frac{1.303206}{2.303206} = 0.5658226.$$

But, what about the variance? Well, if we look at the β estimates, **MARK** reports that the standard error for the estimate of β corresponding to survival is 0.1446688. If we simply back-transform this from the logit scale to the probability scale, we get:

$$\widehat{\text{SE}} = \frac{e^{0.1446688}}{1 + e^{0.1446688}} = \frac{1.155657}{2.155657} = 0.5361043.$$

However, **MARK** reports the estimated standard error for φ as 0.0355404, which isn’t even remotely close to our back-transformed value of 0.5361043.

What has happened? Well, hopefully you now realize that you’re ‘transforming’ the estimate from one scale (logit) to another (probability). And, since you’re working with a ‘transformation’, you need to use the Delta method to estimate the variance of the back-transformed parameter.

Since

$$\hat{\varphi} = \frac{e^{\hat{\beta}}}{1 + e^{\hat{\beta}}},$$

then

$$\widehat{\text{Var}}(\hat{\varphi}) \approx \left(\frac{\partial \hat{\varphi}}{\partial \hat{\beta}} \right)^2 \times \widehat{\text{Var}}(\hat{\beta})$$

$$\begin{aligned}
&= \left(\frac{e^{\hat{\beta}}}{1 + e^{\hat{\beta}}} - \frac{(e^{\hat{\beta}})^2}{1 + (e^{\hat{\beta}})^2} \right)^2 \times \widehat{\text{Var}}(\hat{\beta}) \\
&= \left(\frac{e^{\hat{\beta}}}{(1 + e^{\hat{\beta}})^2} \right)^2 \times \widehat{\text{Var}}(\hat{\beta}).
\end{aligned}$$

It is again worth noting that if

$$\hat{\varphi} = \frac{e^{\hat{\beta}}}{1 + e^{\hat{\beta}}},$$

then it can be easily shown that

$$\hat{\varphi}(1 - \hat{\varphi}) = \frac{e^{\hat{\beta}}}{(1 + e^{\hat{\beta}})^2},$$

which is the derivative of φ with respect to β .

So, we could rewrite our expression for the variance of $\hat{\varphi}$ conveniently as

$$\begin{aligned}
\widehat{\text{Var}}(\hat{\varphi}) &\approx \left(\frac{e^{\hat{\beta}}}{(1 + e^{\hat{\beta}})^2} \right)^2 \times \widehat{\text{Var}}(\hat{\beta}) \\
&= \left(\hat{\varphi}(1 - \hat{\varphi}) \right)^2 \times \widehat{\text{Var}}(\hat{\beta}).
\end{aligned}$$

From **MARK**, the estimate of the SE for $\hat{\beta}$ was 0.1446688. Thus, the estimate of $\text{Var}(\beta)$ is $(0.1446688)^2 = 0.02092906$. Given the estimate of $\hat{\beta}$ of 0.2648275, we substitute into the preceding expression, which yields

$$\begin{aligned}
\widehat{\text{Var}}(\hat{\varphi}) &\approx \left(\frac{e^{\hat{\beta}}}{(1 + e^{\hat{\beta}})^2} \right)^2 \times \widehat{\text{Var}}(\hat{\beta}) \\
&= (0.0603525 \times 0.02092906) \\
&= 0.001263.
\end{aligned}$$

So, the estimated SE for $\hat{\varphi}$ is $\sqrt{0.001263} = 0.0355404$, which is what is reported by **MARK**.

begin sidebar

SE and 95% CI

The standard approach to calculating 95% confidence limits for some parameter θ is $\theta \pm (1.96 \times \text{SE})$. Is this how **MARK** calculates the 95% CI on the real probability scale? Take the example we just considered – the estimated SE for $\hat{\varphi} = 0.5658226$ was $\sqrt{0.001263} = 0.0355404$. So, you might assume that the 95% CI on the real probability scale would be $0.5658226 \pm (2 \times 0.0355404)$: [0.4947418, 0.6369034].

However, this is not what is reported by **MARK** – [0.4953193, 0.6337593], which is quite close, but not exactly the same. Why the difference? The difference is because **MARK** first calculated the 95% CI on the logit scale, before back-transforming to the real probability scale. So, for our estimate of $\hat{\varphi}$, the 95% CI on the logit scale for $\hat{\beta} = 0.2648275$ is $[-0.0187234, 0.5483785]$, which, when back-transformed

to the real probability scale is [0.4953193, 0.6337593], which is what is reported by **MARK**.

In this case, the very small difference between the two CI's is because the parameter estimate was quite close to 0.5. In such cases, not only will the 95% CI be nearly the same (for estimates of 0.5, it will be identical), but they will also be symmetrical.

However, because the logit transform is not linear, the *reconstituted* 95% CI will not be symmetrical around the parameter estimate, especially for parameters estimated near the [0, 1] boundaries. For example, consider the estimate for $\hat{p} = 0.9231757$. On the logit scale, the 95% CI for the β corresponding to p ($\widehat{SE} = 0.5120845$) is [1.4826128, 3.4899840]. The back-transformed CI is [0.8149669, 0.9704014], which is what is reported by **MARK**. This CI is clearly **not** symmetric around $\hat{p} = 0.9231757$. The degree of asymmetry is a function of how close the estimated parameter is to either the 0 or 1 boundary.

Further, the estimated variance for \hat{p} :

$$\begin{aligned}\widehat{\text{Var}}(\hat{p}) &\approx [\hat{p}(1 - \hat{p})]^2 \times \widehat{\text{Var}}(\hat{\beta}) \\ &= [0.9231757(1 - 0.9231757)]^2 \times 0.262231 \\ &= 0.001319,\end{aligned}$$

yields an estimated SE of 0.036318 on the normal probability scale (which is what is reported by **MARK**).

Estimating the 95% CI on the probability scale as $0.9231757 \pm (2 \times 0.036318)$ yields [0.85054, 0.99581], which is clearly quite a bit different, and more symmetrical, than what is reported by **MARK** (from above, [0.8149669, 0.9704014]). **MARK** uses the back-transformed CI to ensure that the reported CI is bounded [0, 1]. As the estimated parameter approaches either the 0 or 1 boundary, the degree of asymmetry in the back-transformed 95% CI that **MARK** reports will increase.

end sidebar

Example (4) – variance of back-transformed estimates - harder

In Chapter 6 we considered the analysis of variation in the survival of the European Dipper, as a function of whether or not there was a flood in the sampling area. Here, we consider just the male Dipper data (the encounter data are contained in `ed_males.inp`). Recall that a flood occurred during over the second and third intervals. For convenience, we'll assume that encounter probability is constant over time, and that survival is a linear function of 'flood'.

Using a logit link function, where 'flood' years were coded in the design matrix using a '1', and 'non-flood' years were coded using a '0', the estimated linear model for survival on the logit scale was:

$$\text{logit}(\hat{\varphi}) = 0.4267863 - 0.5066372(\text{flood})$$

So, in a flood year:

$$\begin{aligned}\text{logit}(\hat{\varphi}_{\text{flood}}) &= 0.4267863 - 0.5066372(\text{flood}) \\ &= 0.4267863 - 0.5066372(1) \\ &= -0.0798509.\end{aligned}$$

Back-transforming onto the real probability scale yields the precise value reported by **MARK**:

$$\hat{\varphi}_{\text{flood}} = \frac{e^{-0.0798509}}{1 + e^{-0.0798509}} = 0.48005.$$

Now, what about the estimated variance for φ_{flood} ? First, what is our 'transformation function' (Y)?

Simple – it is the ‘back-transform’ of the linear equation on the logit scale.

Given that:

$$\begin{aligned}\text{logit}(\hat{\phi}) &= \beta_1 + \beta_2(\text{flood}) \\ &= 0.4267863 - 0.5066372(\text{flood}),\end{aligned}$$

then the back-transform function Y is

$$\hat{Y} = \frac{e^{0.4267863 - 0.5066372(\text{flood})}}{1 + e^{0.4267863 - 0.5066372(\text{flood})}}.$$

Second, since our transformation clearly involves multiple parameters (β_1, β_2) , the estimate of the variance is given to first-order by

$$\begin{aligned}\widehat{\text{Var}}(\hat{Y}) &\approx \mathbf{D}\hat{\Sigma}\mathbf{D}^\top \\ &= \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right] \cdot \widehat{\Sigma} \cdot \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right]^\top.\end{aligned}$$

Given our linear (transformation) equation, then the vector of partial derivatives is (we’ve transposed it to make it easily fit on the page):

$$\begin{aligned}&\left[\left(\frac{\partial(\hat{Y})}{\partial\hat{\beta}_1} \right) \quad \left(\frac{\partial(\hat{Y})}{\partial\hat{\beta}_2} \right) \right]^\top \\ &= \left[\begin{array}{c} \frac{e^{\beta_1 + \beta_2(\text{flood})}}{1 + e^{\beta_1 + \beta_2(\text{flood})}} - \frac{\left(e^{\beta_1 + \beta_2(\text{flood})} \right)^2}{\left(1 + e^{\beta_1 + \beta_2(\text{flood})} \right)^2} \\ \frac{\text{flood} \times e^{\beta_1 + \beta_2(\text{flood})}}{1 + e^{\beta_1 + \beta_2(\text{flood})}} - \frac{\text{flood} \times \left(e^{\beta_1 + \beta_2(\text{flood})} \right)^2}{\left(1 + e^{\beta_1 + \beta_2(\text{flood})} \right)^2} \end{array} \right]\end{aligned}$$

While this is fairly ‘ugly’ looking, the structure is quite straightforward – the only difference between the 2 elements of the vector is that the numerator of both terms (on either side of the minus sign) are multiplied by 1, and flood , respectively. Where do these scalar multipliers come from? They’re simply the partial derivatives of the linear model (we’ll call it Y) on the logit scale:

$$Y = \text{logit}(\hat{\phi}) = \beta_1 + \beta_2(\text{flood}),$$

with respect to each of the parameters (β_i) in turn. In other words, $\partial Y / \partial \beta_1 = 1$, and $\partial Y / \partial \beta_2 = \text{flood}$.

Substituting in our estimates for $\hat{\beta}_1 = 0.4267863$ and $\hat{\beta}_2 = -0.5066372$, and setting $\text{flood}=1$ (to indicate a ‘flood year’) yields:

$$\left[\left(\frac{\partial(\hat{Y})}{\partial\hat{\beta}_1} \right) \quad \left(\frac{\partial(\hat{Y})}{\partial\hat{\beta}_2} \right) \right] = [0.249602 \quad 0.249602]$$

From the **MARK** output (after exporting to a dBase file – and not to the Notepad – in order to get full precision), the full V-C matrix for the parameters β_1 and β_2 is:

$$\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) = \widehat{\Sigma} = \begin{bmatrix} 0.0321405326 & -0.0321581167 \\ -0.0321581167 & 0.0975720877 \end{bmatrix}.$$

So,

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \begin{bmatrix} 0.249602 & 0.249602 \end{bmatrix} \times \begin{bmatrix} 0.0321405326 & -0.0321581167 \\ -0.0321581167 & 0.0975720877 \end{bmatrix} \times \begin{bmatrix} 0.249602 \\ 0.249602 \end{bmatrix} \\ &= 0.0040742678. \end{aligned}$$

The estimated SE for the variance for the reconstituted value of survival for an individual during a ‘flood year’ is $\sqrt{0.0040742678} = 0.0638300$, which is what is reported by **MARK** (to within rounding error).

begin sidebar

Once again...SE and 95% CI

As noted in the preceding example, the standard approach to calculating 95% confidence limits for some parameter θ is $\theta \pm (1.96 \times \text{SE})$. However, to guarantee that the calculated 95% CI is $[0, 1]$ bounded for parameters that are $[0, 1]$ bounded (like ϕ), **MARK** first calculates the 95% CI on the logit scale, before back-transforming to the real probability scale. However, because the logit transform is not linear, the *reconstituted* 95% CI will not be symmetrical around the parameter estimate, especially for parameters estimated near the $[0, 1]$ boundaries.

For the present example, the estimated value of survival for an individual during a ‘flood year’ ($\hat{\phi}_{\text{flood}} = 0.48005$), **MARK** reports a 95% CI of $[0.3586850, 0.6038121]$. But, where do the values $[0.3586850, 0.6038121]$ come from? Clearly, they are not based on $0.48005 \pm 1.96(\text{SE})$. Given $\widehat{\text{SE}} = 0.06383$, this would yield a 95% CI of $[0.35494, 0.60516]$, which is close, but not exactly what **MARK** reports.

In order to derive the 95% CI, we first need to calculate the variance (and SE) of the estimate *on the logit scale*. In the preceding example, this was very straightforward, since the model we considered had a single β term for the parameter of interest. Meaning, we could simply use the estimated SE for β to derive the 95% CI on the logit scale, which we then back-transformed onto the real probability scale.

For the present example, however, the parameter is estimated from a function (transformation) involving more than one β term. In this example, the linear equation, which for consistency with the preceding we will denote as Y , was written as:

$$\hat{Y} = \text{logit}(\hat{\phi}) = \beta_1 + \beta_2(\text{flood})$$

Thus, the estimated variance of $\text{logit}(\hat{\phi}_{\text{flood}})$ is approximated as

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \mathbf{D} \widehat{\Sigma} \mathbf{D}^T \\ &= \begin{bmatrix} \frac{\partial(\hat{Y})}{\partial(\hat{\beta}_1)} & \frac{\partial(\hat{Y})}{\partial(\hat{\beta}_2)} \end{bmatrix} \cdot \widehat{\Sigma} \cdot \begin{bmatrix} \frac{\partial(\hat{Y})}{\partial(\hat{\beta}_1)} & \frac{\partial(\hat{Y})}{\partial(\hat{\beta}_2)} \end{bmatrix}^T. \end{aligned}$$

Since

$$\begin{bmatrix} \frac{\partial(\hat{Y})}{\partial(\hat{\beta}_1)} & \frac{\partial(\hat{Y})}{\partial(\hat{\beta}_2)} \end{bmatrix} = [1 \quad \text{flood}] = [1 \quad 1],$$

and the VC matrix for $\hat{\beta}_1$ and $\hat{\beta}_2$ is

$$\widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) = \widehat{\Sigma} = \begin{bmatrix} 0.0321405356 & -0.0321581194 \\ -0.0321581194 & 0.0975720908 \end{bmatrix},$$

then

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \mathbf{D}\widehat{\Sigma}\mathbf{D}^\top \\ &= \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0.0321405356 & -0.0321581194 \\ -0.0321581194 & 0.0975720908 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &= 0.065396. \end{aligned}$$

So, the $\widehat{\text{SE}}$ – on the logit scale! – is $\sqrt{0.065396} = 0.255727$. Thus, the 95% CI on the estimate on the logit scale, $\text{logit}(\hat{\phi}_{\text{flood}}) = -0.0798509 \pm 1.96(0.255727) = [-0.581076, 0.421374]$.

All that is left is to back-transform the limits on the CI to the real probability scale:

$$[-0.94321, 0.78351] \longrightarrow \left[\frac{e^{-0.581076}}{1 + e^{-0.581076}}, \frac{e^{0.421374}}{1 + e^{0.421374}} \right] = [0.358685, 0.603812],$$

which is what is reported by **MARK** (to within rounding error).

end sidebar

Example (5) – variance of back-transformed estimates – harder still

In Chapter 11, we considered analysis of the effect of various functions of mass, m , and mass-squared, m^2 , on the survival of a hypothetical species of bird (the simulated data are in file **indcov1.inp**). The linear function relating survival to m and m^2 , on the logit scale, is:

$$\text{logit}(\hat{\phi}) = 0.2567341 + 1.1750463(m_s) - 1.0554957(m_s^2).$$

Note that for the two mass terms in the equation, there is a small subscript ‘s’, reflecting the fact that these are ‘standardized’ masses. Recall that we standardized the covariates by subtracting the mean of the covariate, and dividing by the standard deviation (the use of standardized or non-standardized covariates is discussed at length in Chapter 11).

Thus, for each individual in the sample, the estimated survival probability (on the logit scale) for that individual, given its mass, is given by:

$$\text{logit}(\hat{\phi}) = 0.2567333 + 1.1750526 \left(\frac{m - \bar{m}}{\text{SD}_m} \right) - 1.0555024 \left(\frac{m^2 - \bar{m}^2}{\text{SD}_{m^2}} \right).$$

In this expression, m refers to mass and m^2 refers to mass2. The output from **MARK** (preceding page) actually gives you the mean and standard deviations for both covariates: for mass, mean = 109.9680, and SD = 24.7926, while for mass2, the mean = 12,707.4640, and the SD = 5,532.0322. The ‘value’ column shows the standardized values for mass and mass2 (0.803 and 0.752) for the first individual in the data file.

Now let’s consider a worked example of the calculation of the variance of estimated survival. Suppose

the mass of the bird was 110 units, so that $m = 110$, $m_2 = (110)^2 = 12,100$.

$$\begin{aligned}\text{logit}(\hat{\varphi}) &= 0.2567333 + 1.1750526 \left(\frac{(110 - 109.9680)}{24.7926} \right) - 1.0555024 \left(\frac{(12,100 - 12,707.4640)}{5,532.0322} \right) \\ &= 0.3742.\end{aligned}$$

So, if $\text{logit}(\hat{\varphi}) = 0.374$, then the reconstituted estimate of φ , transformed back from the logit scale is:

$$\frac{e^{0.374152}}{1 + e^{0.374152}} = 0.5925.$$

Thus, for an individual weighing 110 units, the expected annual survival probability is approximately 0.5925 (which is what **MARK** reports if you use the '**User specify covariate**' option).

What about the variance (and corresponding SE) for this estimate? First, what is our 'transformation function' (Y)? For the present example, it is the 'back-transform' of the linear equation on the logit scale. Given that:

$$\begin{aligned}\text{logit}(\hat{\varphi}) &= \beta_1 + \beta_2(m_s) + \beta_3(m_s^2) \\ &= 0.2567333 + 1.1750526(m_s) - 1.0555024(m_s^2),\end{aligned}$$

then the back-transform Y is:

$$\hat{Y} = \frac{e^{0.2567333 + 1.1750526(m_s) - 1.0555024(m_s^2)}}{1 + e^{0.2567333 + 1.1750526(m_s) - 1.0555024(m_s^2)}}$$

As in the preceding example, since our transformation clearly involves multiple parameters ($\beta_1, \beta_2, \beta_3$), the estimate of the variance is given by:

$$\begin{aligned}\widehat{\text{Var}}(\hat{Y}) &\approx \mathbf{D}\widehat{\Sigma}\mathbf{D}^T \\ &= \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right] \cdot \widehat{\Sigma} \cdot \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right]^T\end{aligned}$$

Given our linear (transformation) equation (from above) then the vector of partial derivatives is:

$$\begin{bmatrix} \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_0} \right) \\ \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_1} \right) \\ \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_2} \right) \end{bmatrix} = \begin{bmatrix} \frac{e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}}{1 + e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}} - \frac{[e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}]^2}{[1 + e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}]^2} \\ \frac{m \times e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}}{1 + e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}} - \frac{m \times [e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}]^2}{[1 + e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}]^2} \\ \frac{m_2 \times e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}}{1 + e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}} - \frac{m_2 \times [e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}]^2}{[1 + e^{\beta_1 + \beta_2(m) + \beta_3(m_2)}]^2} \end{bmatrix}.$$

Although this looks complicated, the structure is actually quite straightforward – the only difference between the 3 elements of the vector is that the numerator of both terms (on either side of the minus

sign) are multiplied by 1, m , and $m2$, respectively, which are simply the partial derivatives of the linear model (we'll call it Y) on the logit scale:

$$\hat{Y} = \text{logit}(\hat{\phi}) = \beta_1 + \beta_2(m_s) + \beta_3(m_s^2),$$

with respect to each of the parameters (β_i) in turn. In other words, $\partial Y / \partial \beta_1 = 1$, $\partial Y / \partial \beta_2 = m$, and $\partial Y / \partial \beta_3 = m2$.

So, now that we have our vectors of partial derivatives of the transformation function with respect to each of the parameters, we can simplify things considerably by substituting in the standardized values for m and $m2$, and the estimated parameter values ($\hat{\beta}_1$, $\hat{\beta}_2$, and $\hat{\beta}_3$).

For a mass of 110 g, the standardized values for m and $m2$ are:

$$m_s = \left(\frac{110 - 109.9680}{24.7926} \right) = 0.0012895,$$

$$m2_s = \left(\frac{12100 - 12707.4640}{5532.0322} \right) = -0.1098085.$$

The estimates for $\hat{\beta}_i$ we read directly from **MARK**:

$$\hat{\beta}_1 = 0.2567333, \hat{\beta}_2 = 1.1750526, \text{ and } \hat{\beta}_3 = -1.0555024.$$

Substituting in these estimates for $\hat{\beta}_i$ and the standardized m and $m2$ values into our vector of partial derivatives (which we've transposed in the following to save space) yields:

$$\left[\left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_1} \right) \quad \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_2} \right) \quad \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_3} \right) \right]^T = \begin{bmatrix} 0.241451 \\ 0.000311 \\ -0.026513 \end{bmatrix}.$$

From the **MARK** output (after exporting to a dBase file – and not to the editor – in order to get full precision), the full V-C matrix for the β parameters is

$$\begin{bmatrix} 0.0009006967 & -0.0004110129 & 0.0003662771 \\ -0.0004110129 & 0.0373928740 & -0.0364291221 \\ 0.0003662771 & -0.0364291221 & 0.0362817338 \end{bmatrix}.$$

So,

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \begin{bmatrix} 0.241451 & 0.000311 & -0.026513 \end{bmatrix} \\ &\times \begin{bmatrix} 0.0009006967 & -0.0004110129 & 0.0003662771 \\ -0.0004110129 & 0.0373928740 & -0.0364291221 \\ 0.0003662771 & -0.0364291221 & 0.0362817338 \end{bmatrix} \times \begin{bmatrix} 0.241451 \\ 0.000311 \\ -0.026513 \end{bmatrix} \\ &= 0.000073867. \end{aligned}$$

So, the estimated SE for Var for the reconstituted value of survival for an individual weighing 110 g is $\sqrt{0.000073867} = 0.0085946$, which is exactly what is reported by **MARK**.

It is important to remember that the estimated variance will vary depending on the mass you use – the estimate of the variance for a 110 g individual (0.000073867) will differ from the estimated variance

for a (say) 120 g individual. For a 120 g individual, the standardized values of m and m_2 are 0.404636 and 0.3059512, respectively.

Based on these values, then:

$$\left[\left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_1} \right) \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_2} \right) \left(\frac{\partial(\hat{Y})}{\partial \hat{\beta}_3} \right) \right]^T = \begin{bmatrix} 0.239817 \\ 0.097039 \\ 0.073372 \end{bmatrix}.$$

Given the variance covariance-matrix for this model (shown above), then

$$\widehat{\text{Var}}(\hat{Y}) \approx \mathbf{D}\Sigma\mathbf{D}^T = 0.000074246.$$

Thus, the estimated SE for the reconstituted value of survival for an individual weighing 120 g is $\sqrt{0.000074246} = 0.0086166$, which again is *exactly* what is reported by **MARK**.

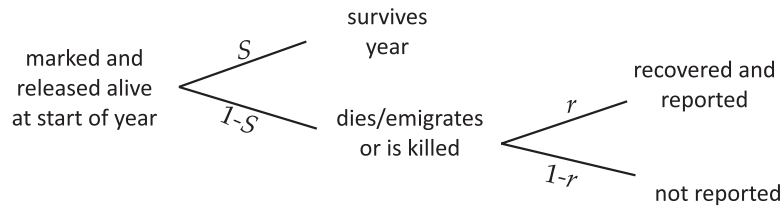
Note that this value for the SE for a 120 g individual (0.008617) differs from the SE estimated for a 110 g individual (0.008595), albeit not by much (the small difference here is because this is a very large simulated data set based on a deterministic model – see Chapter 11 for details). Since each weight would have its own estimated survival, and associated estimated variance and SE, to generate a curve showing the reconstituted values and their SE, you'd need to iteratively calculate $\mathbf{D}\Sigma\mathbf{D}^T$ over a range of weights. We'll leave it to you to figure out how to handle the programming if you want to do this on your own. For the less ambitious, **MARK** has the capacity to do much of this for you – you can output the 95% CI 'data' over a range of individual covariate values to a spreadsheet (see section 11.5 in Chapter 11).

Example (6) – estimating variance + covariance in transformations

Here, we consider application of the Delta method to joint estimation of the variance of a parameter, and the *covariance* of that parameter with another, where one of the two parameters is a linear transformation of the other. This is a somewhat complicated but quite useful example, since it illustrates how you can use the Delta method to estimate not only the variance of individual parameters, but the covariance structure among parameters as well.

There are many instances where the magnitude of the covariance among parameters is of particular interest. Here, we consider such a situation, in terms of different parameterizations for analysis of dead recovery data. Dead recovery models are covered in detail in Chapter 8 – here, we briefly review two different parameterizations (the 'Seber' and 'Brownie' parameterizations), and the context of our interest in the covariance between two different parameters.

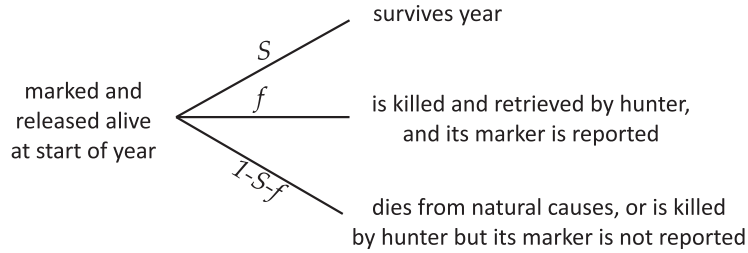
The encounter process for the Seber parameterization (1973: 254) is illustrated in the following:



Marked individuals are assumed to survive from release (i) to ($i + 1$) with probability S_i . Individuals

may die during the interval, either due to harvest or to 'natural' mortality. The probability that dead marked individuals are reported during each period (i) between releases, and (most generally) where the death is not necessarily related to harvest, is r_i . In other words, r_i is the joint probability of (i) the marked individual dying from either harvest or natural causes, and (ii) being recovered and reported.

Brownie *et al.* (1985) (hereafter, simply 'Brownie') developed a different parameterization for dead recovery data, where the sources of mortality (harvest, versus 'natural' or non-harvest) are modeled separately. The encounter process for this parameterization is shown below:



In the Brownie parameterization, S_i is the probability that the individual survives the interval from release occasion (i) to ($i + 1$) (note that the definition for the probability of survival is logically identical between the Seber and Brownie parameterizations). The probability that the individual dies from either source of mortality is simply $(1 - S)$.

However, in contrast to the Seber parameterization, Brownie specified a parameter f , to represent the probability that an individual dies specifically due to harvest during interval (i), and is reported ('encountered'). Thus, the probability that the individuals dies from natural causes is $(1 - S - f)$.

Under the Seber parameterization, the probability of the encounter history '11' is $r(1 - S)$. Under the Brownie parameterization, the expected probability of this event is simply f . Since the encounter history is the same, we can set the different parameterizations for the expected probability of the event equal to each other, generating the following expressions relating the two parameterizations:

$$f_i = r_i(1 - S_i) \quad r_i = \frac{f_i}{(1 - S_i)}$$

Clearly, the parameter r_i is a reduced parameter, and can be expressed as a function of two other parameters normally found in the Brownie parameterization. An obvious practical question is, why choose one parameterization over the other, and does it matter?

This issue is discussed more fully in Chapter 8, but for now, we focus on the left-hand expression:

$$f_i = r_i(1 - S_i).$$

So, given estimates of \hat{r}_i and \hat{S}_i from a Seber analysis, we could use this algebraic relationship (i.e., transformation) and the Delta method to generate estimates of \hat{f} and estimate $\widehat{\text{Var}}(\hat{f})$.

However, in addition, suppose we are also interested in estimating the covariance $\widehat{\text{Cov}}(\hat{f}, \hat{S})$. Recall from above that the parameter f relates in part to the probability of being harvested. We might naturally be interested in the relationship between harvest mortality f , and overall annual survival, S . For example, if harvest and natural mortality are strictly additive, then we might expect a negative covariance between survival and harvest (i.e., as the probability of mortality due to harvest increases, annual survival will decrease). Whether or not the covariance is negative has important implications for harvest management (see full discussion in the Williams, Nichols & Conroy 2001 book).

We'll begin by considering estimation of the variance for \hat{f} only, using the Delta method. Let the transformation g be $f = (1 - S)r$.

Given \hat{S} , \hat{r} , $\widehat{\text{Var}}(\hat{S})$ and $\widehat{\text{Var}}(\hat{r})$, then the Jacobian for g is

$$\begin{bmatrix} \frac{\partial g}{\partial S} & \frac{\partial g}{\partial r} \end{bmatrix} = \begin{bmatrix} -\hat{r} & 1 - \hat{S} \end{bmatrix},$$

and thus

$$\widehat{\text{Var}}(\hat{f}) \approx \begin{bmatrix} -\hat{r} & 1 - \hat{S} \end{bmatrix} \cdot \widehat{\Sigma} \cdot \begin{bmatrix} -\hat{r} \\ 1 - \hat{S} \end{bmatrix},$$

where $\widehat{\Sigma}$ is the variance-covariance matrix for S and r :

$$\widehat{\Sigma} = \begin{bmatrix} \widehat{\text{Var}}(\hat{S}) & \widehat{\text{Cov}}(\hat{S}, \hat{r}) \\ \widehat{\text{Cov}}(\hat{S}, \hat{r}) & \widehat{\text{Var}}(\hat{r}) \end{bmatrix}.$$

So,

$$\widehat{\text{Var}}(\hat{f}) \approx \begin{bmatrix} -\hat{r} & 1 - \hat{S} \end{bmatrix} \cdot \widehat{\Sigma} \cdot \begin{bmatrix} -\hat{r} \\ 1 - \hat{S} \end{bmatrix},$$

yielding

$$\widehat{\text{Var}}(\hat{f}) \approx S^2 \cdot \widehat{\text{Var}}(\hat{r}) + 2\hat{S} \cdot \hat{r} \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) + \hat{r}^2 \cdot \widehat{\text{Var}}(\hat{S}) - 2\hat{S} \cdot \widehat{\text{Var}}(\hat{r}) - 2\hat{r} \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) + \widehat{\text{Var}}(\hat{r}),$$

which, with a little re-arranging, simplifies to

$$\widehat{\text{Var}}(\hat{f}) \approx \hat{r}^2 \cdot \widehat{\text{Var}}(\hat{S}) - 2[(1 - \hat{S})\hat{r}] \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) + (1 - \hat{S})^2 \cdot \widehat{\text{Var}}(\hat{r}).$$

So, given estimates of $\widehat{\text{Var}}(\hat{S})$, $\widehat{\text{Var}}(\hat{r})$ and $\widehat{\text{Cov}}(\hat{S}, \hat{r})$ from a 'Seber' analysis of dead recovery data, we can derive an estimate of $\widehat{\text{Var}}(\hat{f})$ for the 'Brownie' parameterization by substituting those estimates into this approximation and solving (we demonstrate this empirically using analysis of dead recovery in **MARK** in the next sidebar-).

Now, what if instead of $\widehat{\text{Var}}(\hat{f})$ only, we are also interested in estimating the *covariance* of (say) f and S ? Such a covariance might be of interest since f is a function of S , and there may be interest in the degree to which S varies as a function of f (see above). Thus, we want to apply the Delta method to a function (in this case, the covariance) of two parameters, f and S .

The key step here is recognizing that there are in fact two different functions (or, transformations) involved, which we'll call g_1 and g_2 :

$$g_1 : S \rightarrow S \quad \text{and} \quad g_2 : (1 - S)r \rightarrow f$$

You might be puzzled by $g_1 : S \rightarrow S$. In fact, this represents a null transformation – a direct, non-transformative 1:1 mapping between S under the Seber parameterization and under the Brownie parameterization (since the probability of surviving is, logically, the same under the two parameterizations). This is analogous to generating the estimate for \hat{S}_i under one parametrization by multiplying the same estimate under the other parameterization by the scalar constant 1.

Thus, with two transformations, we generate a Jacobian *matrix* of partial derivatives of each of the transformations with respect to S and r :

$$\begin{aligned} \begin{bmatrix} \frac{\partial g_1}{\partial \hat{S}} & \frac{\partial g_1}{\partial \hat{r}} \\ \frac{\partial g_2}{\partial \hat{S}} & \frac{\partial g_2}{\partial \hat{r}} \end{bmatrix} &= \begin{bmatrix} \frac{\partial \hat{S}}{\partial \hat{S}} & \frac{\partial \hat{S}}{\partial \hat{r}} \\ \frac{\partial \hat{f}}{\partial \hat{S}} & \frac{\partial \hat{f}}{\partial \hat{r}} \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ -r & 1 - \hat{S} \end{bmatrix}. \end{aligned}$$

Given the variance-covariance matrix $\widehat{\Sigma}_{\hat{S}, \hat{r}}$ for \hat{S} and \hat{r}

$$\widehat{\Sigma}_{\hat{S}, \hat{r}} = \begin{bmatrix} \widehat{\text{Var}}(\hat{S}) & \widehat{\text{Cov}}(\hat{S}, \hat{r}) \\ \widehat{\text{Cov}}(\hat{S}, \hat{r}) & \widehat{\text{Var}}(\hat{r}) \end{bmatrix},$$

we evaluate the sampling variance-covariance matrix for \hat{S} and \hat{f} as the matrix product

$$\widehat{\Sigma}_{\hat{S}, \hat{f}} = \begin{bmatrix} 1 & 0 \\ -r & 1 - \hat{S} \end{bmatrix} \cdot \widehat{\Sigma}_{\hat{S}, \hat{r}} \cdot \begin{bmatrix} 1 & -r \\ 0 & 1 - \hat{S} \end{bmatrix},$$

which (after a bit of algebra) yields

$$\widehat{\Sigma}_{\hat{S}, \hat{f}} = \begin{bmatrix} \widehat{\text{Var}}(\hat{S}) & -r \cdot \widehat{\text{Var}}(\hat{S}) + (1 - \hat{S}) \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) \\ -r \cdot \widehat{\text{Var}}(\hat{S}) + (1 - \hat{S}) \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) & \hat{r}^2 \cdot \widehat{\text{Var}}(\hat{S}) - 2[(1 - \hat{S})\hat{r}] \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) + (1 - \hat{S})^2 \cdot \widehat{\text{Var}}(\hat{r}) \end{bmatrix}.$$

The diagonal matrix elements [1, 1] and [2, 2] are the expressions for the approximate variance of \hat{S} and \hat{f} , respectively (note that the expression in element [2, 2], for $\widehat{\text{Var}}(\hat{f})$, is identical to the expression we derived on the preceding page). Off-diagonal elements [1, 2] and [2, 1] (which are the same) are the expressions for the approximate *covariance* of \hat{f} and \hat{S} .

As noted earlier, interpretation of the estimated variance and covariance is dependent on the source of the variance-covariance matrix, $\widehat{\Sigma}$, used in the calculations. If $\widehat{\Sigma}$ is constructed using variances and covariances from the usual ML parameter estimates, then the resulting estimate for variance is an estimate of the *total* variance (i.e., sampling + process, where process variation represents the underlying ‘biological’ variation). In contrast, if $\widehat{\Sigma}$ is based on estimated *process* variances and covariances only, then the estimate for variance is an estimate of the *process* variance. Decomposition of total variance into sampling and process components is covered in detail in Appendix D.

begin sidebar

‘proof’ by empirical example...

In the preceding, we used the analytical approach and the Delta method to derive expressions for variances and covariances. Somewhat elegant algebra, but, does it actually work?

Here use **MARK** to analyze a simulated dead recovery data set (**brownie_seber.inp**): 7 occasions, 2,000 individuals marked and released at each occasion, simulated under the true generating model $\{S, p\}$ (i.e., no temporal variation in either parameter) using the ‘Seber’ dead recovery data type. The

parameter values used in the simulation were $S = 0.8$, $r = 0.4$.

We start by fitting the ‘true’ generating model $\{S, p.\}$ to these data, under the Seber data type. Here are the parameter estimates (quite close to the true values used in simulating the data):

Real Function Parameters of $\{S(.)r(.)$ - Seber}				
Parameter	Estimate	Standard Error	95% Confidence Lower	Interval Upper
1:S	0.8001978	0.0094206	0.7810923	0.8180243
2:r	0.4005510	0.0135221	0.3743616	0.4273211

Here is the corresponding variance-covariance matrix output by **MARK**, with $\widehat{\text{Var}}(\hat{S})$ and $\widehat{\text{Var}}(\hat{r})$ along the diagonal, respectively, and $\widehat{\text{Cov}}(\hat{S}, \hat{r})$ on the off-diagonal:

```

0.0000887478    0.0001129638
0.0001129638    0.0001828483

```

Earlier, we used the Delta method to derive an estimate of $\widehat{\text{Var}}(\hat{f})$ as:

$$\widehat{\text{Var}}(\hat{f}) \approx \hat{r}^2 \cdot \widehat{\text{Var}}(\hat{S}) - 2[(1 - \hat{S})\hat{r}] \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) + (1 - \hat{S})^2 \cdot \widehat{\text{Var}}(\hat{r}).$$

Substituting in the parameter estimates, we can derive an estimate for the variance of \hat{f} :

$$\begin{aligned}
\widehat{\text{Var}}(\hat{f}) &\approx \hat{r}^2 \cdot \widehat{\text{Var}}(\hat{S}) - 2[(1 - \hat{S})\hat{r}] \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) + (1 - \hat{S})^2 \cdot \widehat{\text{Var}}(\hat{r}) \\
&= [(0.4005510)^2 \cdot (0.0000887478)] - [2[(1 - 0.8001978) \cdot 0.4005510] \cdot 0.0001129638] \\
&\quad + [(1 - 0.8001978)^2 \cdot 0.0001828483] \\
&= 0.0000034571.
\end{aligned}$$

Next, in the same **MARK** session, we select ‘**PIM | change data type**’ and switch from the ‘Seber’ data type → the ‘Brownie’ data type. In other words, we’re now going to model the same recovery data, in the same **MARK** session, but this time using the Brownie parameterization. We fit the Brownie model $\{S, f.\}$ to the data. The model deviances and AIC values will be identical to the model fit under the Seber parameterization – our interest is in the variance-covariance matrix. We see in the following

```

0.0000887478    -0.0000129776
-0.0000129776    0.0000034571

```

that the estimate $\widehat{\text{Var}}(\hat{f}) = 0.0000034571$ from fitting the Brownie model (the [2,2] element, above) is identical to the value approximated using the Delta method (0.0000034571).

We also derived an approximate of the covariance of \hat{S} and \hat{f} as

$$\begin{aligned}
\widehat{\text{Cov}}(\hat{S}, \hat{f}) &\approx -r \cdot \widehat{\text{Var}}(\hat{S}) + (1 - \hat{S}) \cdot \widehat{\text{Cov}}(\hat{S}, \hat{r}) \\
&= -(0.4005510)(0.0000887478) + [(1 - 0.8001978) \cdot 0.0001129638] \\
&= -0.00001297760428
\end{aligned}$$

which again is exactly the same as the value from the Brownie analysis (the [1,2] or equivalent [2,1] element, above).

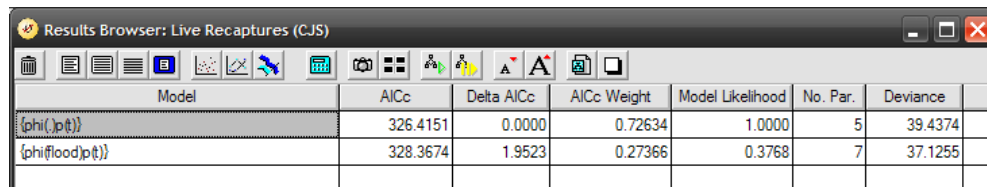
It is perhaps fair to note that if all we wanted was an estimate of the covariance of \hat{S} and \hat{f} , we could have simply analyzed the data using the ‘Brownie’ data type in the first place. While this is true for this particular question, the larger point we’ve just demonstrated is that we can apply the Delta method in a fairly straightforward way to estimate not just the variance for a derived parameter, but the covariances among them as well.

end sidebar

B.6. Delta method and model averaging

In the preceding examples, we focused on the application of the Delta method to transformations of parameter estimates from a single model. However, as introduced in Chapter 4 – and emphasized throughout the remainder of this book – we’re often interested in accounting for model selection uncertainty by using model-averaged values. There are no major complications for application of the Delta method to model-averaged parameter values – you simply need to make sure you use model-averaged values for each element of the calculations.

We’ll demonstrate this using analysis of the male dipper data (`ed_males.inp`). Suppose that we fit 2 candidate models to these data: $\{\varphi, p_t\}$ and $\{\varphi_{flood} p_t\}$. In other words, a model where survival is constant over time, and a model where survival is constrained to be a function of a binary ‘flood’ variable (see section 6.4 of Chapter 6). Here are the results of fitting these 2 models to the data:



Model	AICc	Delta AICc	AICc Weight	Model Likelihood	No. Par.	Deviance
{phi(.), p(t)}	326.4151	0.0000	0.72634	1.0000	5	39.4374
{phi(flood), p(t)}	328.3674	1.9523	0.27366	0.3768	7	37.1255

As expected (based on the analysis of these data presented in Chapter 6), we see that there is some evidence of model selection uncertainty – the model where survival is constant over time has roughly 2-3 times the weight as does the ‘flood’ model.

The model averaged values for each interval are shown below:

	1	2	3	4	5	6
<i>estimate</i>	0.5673	0.5332	0.5332	0.5673	0.5673	0.5673
<i>SE</i>	0.0441	0.0581	0.0581	0.0441	0.0441	0.0441

Now, suppose we want to derive the best estimate of the probability of survival over (say) the first 3 intervals. Clearly, all we need to do is take the product of the 3 model-averaged values corresponding to the first 3 intervals:

$$(0.5673 \times 0.5332 \times 0.5332) = 0.1613.$$

In other words, our best estimate of the probability that a male dipper would survive from the start of the time series to the end of the third interval is 0.1613.

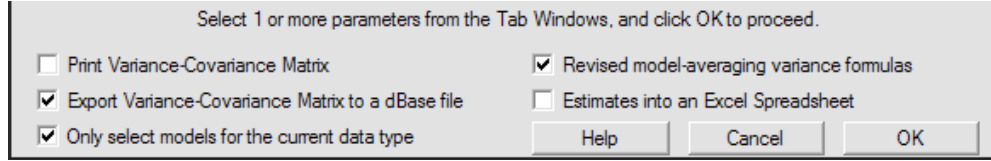
What about the standard error of this product? Here, we use the Delta method. Recall that:

$$\begin{aligned} \widehat{\text{Var}}(\hat{Y}) &\approx \mathbf{D}\widehat{\Sigma}\mathbf{D}^T \\ &= \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right] \cdot \widehat{\Sigma} \cdot \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right]^T, \end{aligned}$$

where Y is some linear or nonlinear function of the parameter estimates $\hat{\theta}_1, \hat{\theta}_2, \dots$. For this example, Y is the product of the survival estimates.

So, the first thing we need to do is to generate the estimated variance-covariance matrix for the model averaged survival estimates. This is easy enough to do – in the ‘**Model Averaging Parameter Selection**’

window, you simply need to ‘**Export Variance-Covariance Matrix to a dBase file**’ - you do this by checking the appropriate check box (lower-left, as shown below):



The ‘rounded’ values which would be output to the Notepad are shown below. (Remember, however, that for the actual calculations, you want to use the full precision variance-covariance matrix from the exported dBase file.)

Unconditional Variance-Covariance Matrix of Model Averaged Estimates Variance-covariance matrix of estimates on diagonal and below, Correlation matrix of estimates above diagonal.			
	1	2	3
1	0.00194	0.04923	0.04923
2	0.00013	0.00337	1.00000
3	0.00013	0.00337	0.00337

All that remains is to substitute our model-averaged estimates for (i) $\hat{\phi}$ and (ii) the variance-covariance matrix (above), into $\widehat{\text{Var}}(\hat{Y}) \approx \mathbf{D}\mathbf{\Sigma}\mathbf{D}^T$.

$$\begin{aligned}
 \widehat{\text{Var}}(\hat{Y}) &\approx \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right] \cdot \widehat{\Sigma} \cdot \left[\frac{\partial(\hat{Y})}{\partial(\hat{\theta})} \right]^T \\
 &= \begin{bmatrix} (\bar{\phi}_2 \bar{\phi}_3) & (\bar{\phi}_1 \bar{\phi}_3) & (\bar{\phi}_1 \bar{\phi}_2) \end{bmatrix} \cdot \begin{bmatrix} \widehat{\text{Var}}(\bar{\phi}_1) & \widehat{\text{Cov}}(\bar{\phi}_1, \bar{\phi}_2) & \widehat{\text{Cov}}(\bar{\phi}_1, \bar{\phi}_3) \\ \widehat{\text{Cov}}(\bar{\phi}_1, \bar{\phi}_2) & \widehat{\text{Var}}(\bar{\phi}_2) & \widehat{\text{Cov}}(\bar{\phi}_2, \bar{\phi}_3) \\ \widehat{\text{Cov}}(\bar{\phi}_3, \bar{\phi}_1) & \widehat{\text{Cov}}(\bar{\phi}_3, \bar{\phi}_2) & \widehat{\text{Var}}(\bar{\phi}_3) \end{bmatrix} \cdot \begin{bmatrix} (\bar{\phi}_2 \bar{\phi}_3) \\ (\bar{\phi}_1 \bar{\phi}_3) \\ (\bar{\phi}_1 \bar{\phi}_2) \end{bmatrix} \\
 &= [0.284303069 \quad 0.3024783390 \quad 0.3024783390] \\
 &\quad \times \begin{bmatrix} 0.0019410083 & 0.0001259569 & 0.0001259569 \\ 0.0001259569 & 0.0033727452 & 0.0033727423 \\ 0.0001259569 & 0.0033727423 & 0.0033727452 \end{bmatrix} \times \begin{bmatrix} 0.284303069 \\ 0.3024783390 \\ 0.3024783390 \end{bmatrix} \\
 &= 0.001435.
 \end{aligned}$$

B.7. Summary

In this appendix, we’ve briefly introduced a convenient, generally straightforward method for deriving an estimate of the sampling variance for transformations of one or more variables. Such transformations are quite commonly encountered when using **MARK**, and having a method to derive estimates of the sampling variances is convenient. The most straightforward method – based on a first-order Taylor

series expansion – is known generally as the ‘Delta method’.

However, as discussed (and demonstrated) a first-order Taylor series approximation may not always be appropriate, especially if the transformation is highly non-linear, and if there is significant variation in the data. In such case, you may have to resort to higher-order approximations, or numerically intensive bootstrapping approaches.

B.8. References

- dos Santos Dias, C. T., Samaranayaka, A., and Manly, B. F. (2008) On the use of correlated beta random variables with animal population modelling. *Ecological Modelling*, **215**, 293-300.
- Ver Hoef, J. M. (2012) Who invented the Delta method? *The American Statistician*, **66**, 124-127.

Addendum B.1 – ‘computationally intensive’ approaches

At the start of this appendix, we motivated the Delta method as an approach for deriving an estimate of the expectation or variance of a function of one or more parameters – specifically, an approach that was *not* ‘compute-intensive’. While this approach has a certain elegance, application to complex functions can be cumbersome. Further, transformations that are strongly nonlinear near the mass of the data may necessitate using a higher-order Taylor series expansion, which again can be complex for a particular function.

In such cases, it is useful to at least be aware of alternative, compute-intensive approaches. Here, we briefly introduce two different approaches, applied to the estimation of the variance of the product of survival estimates, using the dipper example presented in section B.4. Again, we’ll use estimates from model $\{\varphi_i p_i\}$ fit to the male European dipper data set, and again, we’ll suppose we’re interested in the probability of surviving from the start of the first interval to the end of the third interval.

As noted in section B.4, the estimate of this probability is easy enough:

$$\begin{aligned}\hat{Y} &= (\hat{\phi}_1 \times \hat{\phi}_2 \times \hat{\phi}_3) \\ &= (0.6109350 \times 0.458263 \times 0.4960239) = 0.138871.\end{aligned}$$

So, the estimated probability of a male Dipper surviving over the first three intervals is $\sim 14\%$ (again, assuming that our time-dependent survival model is a valid model).

i. using the Delta method...

To derive the estimate of the variance of the product using the Delta method, we require the variance-covariance matrix for the survival estimates:

$$\begin{aligned}\widehat{\text{Cov}}(\hat{Y}) &= \widehat{\Sigma} = \begin{bmatrix} \widehat{\text{Var}}(\hat{\phi}_1) & \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_2) & \widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_3) \\ \widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_1) & \widehat{\text{Var}}(\hat{\phi}_2) & \widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_3) \\ \widehat{\text{Cov}}(\hat{\phi}_3, \hat{\phi}_1) & \widehat{\text{Cov}}(\hat{\phi}_3, \hat{\phi}_2) & \widehat{\text{Var}}(\hat{\phi}_3) \end{bmatrix} \\ &= \begin{bmatrix} 0.0224330125 & -0.0003945405 & 0.0000654469 \\ -0.0003945405 & 0.0099722201 & -0.0002361998 \\ 0.0000654469 & -0.0002361998 & 0.0072418858 \end{bmatrix}.\end{aligned}$$

For this example, the transformation we’re applying to our 3 survival estimates (which we’ll call Y) is the product of the estimates (i.e., $\hat{Y} = \hat{\phi}_1 \hat{\phi}_2 \hat{\phi}_3$).

Thus, our variance estimate is given as

$$\widehat{\text{Var}}(\hat{Y}) \approx \begin{bmatrix} \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_1}\right) & \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_2}\right) & \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_3}\right) \end{bmatrix} \cdot \widehat{\Sigma} \cdot \begin{bmatrix} \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_1}\right) \\ \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_2}\right) \\ \left(\frac{\partial \hat{Y}}{\partial \hat{\phi}_3}\right) \end{bmatrix}.$$

Each of the partial derivatives for \hat{Y} is easy enough to derive for this example. Since $\hat{Y} = \hat{\phi}_1 \hat{\phi}_2 \hat{\phi}_3$, then $\partial \hat{Y} / \partial \hat{\phi}_1 = \hat{\phi}_2 \hat{\phi}_3$. And so on.

Expanding the preceding results in:

$$\begin{aligned}\widehat{\text{Var}}(\hat{Y}) &\approx \hat{\phi}_2^2 \hat{\phi}_3^2 [\widehat{\text{Var}}(\hat{\phi}_1)] + \hat{\phi}_1^2 \hat{\phi}_3^2 [\widehat{\text{Var}}(\hat{\phi}_2)] + \hat{\phi}_1^2 \hat{\phi}_2^2 [\widehat{\text{Var}}(\hat{\phi}_3)] \\ &\quad + 2\hat{\phi}_2 \hat{\phi}_3 \hat{\phi}_1 [\widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_2)] + 2\hat{\phi}_2^2 \hat{\phi}_3 \hat{\phi}_1 [\widehat{\text{Cov}}(\hat{\phi}_1, \hat{\phi}_3)] + 2\hat{\phi}_1^2 \hat{\phi}_3 \hat{\phi}_2 [\widehat{\text{Cov}}(\hat{\phi}_2, \hat{\phi}_3)].\end{aligned}$$

After substituting in our estimates for ϕ_i and the variances and covariances, our estimate for the variance of the product $\hat{Y} = (\hat{\phi}_1 \hat{\phi}_2 \hat{\phi}_3)$ is (to first-order) $\widehat{\text{Var}}(Y) = 0.0025565$.

Now, we consider a couple of ‘compute-intensive’ approaches.

ii. simulation from a multivariate normal distribution...

The basic idea behind the approach we illustrate here is straightforward: we (i) simulate data as random draws from a multivariate normal distribution with known means and variance-covariance, (ii) generate the product of these random draws, and (iii) derive numerical estimates of the mean (expectation) and variance of these products.

However, we need to be a bit careful here. If we simulate the random normal draws on the real probability scale, then we run the risk of simulating random values which are not plausible, because they fall outside the $[0, 1]$ interval (e.g., you could simulate a survival probability > 1 , or < 0 , neither of which are possible). To circumvent this problem, we simulate random normal variables on the logit scale (i.e., logit-normal deviates) using the β estimates and the variance-covariance matrix (both estimated on the logit scale), back-transform the random deviates from the logit \rightarrow real probability scale, and then generate the product on the real probability scale.*

For the male Dipper data, the β estimates generated using an identity design matrix (such that each estimated β corresponds to the survival estimate for that interval – see Chapter 6 for specifics) are: $\hat{\beta}_1 = 0.4512441$, $\hat{\beta}_2 = -0.1673372$, $\hat{\beta}_3 = -0.0159047$. The variance-covariance matrix for the β estimates (which can be output from **MARK**) is:

$$\begin{aligned}\widehat{\text{Cov}}(\hat{Y}) &= \sum \\ &= \begin{bmatrix} \widehat{\text{Var}}(\hat{\beta}_1) & \widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_2) & \widehat{\text{Cov}}(\hat{\beta}_1, \hat{\beta}_3) \\ \widehat{\text{Cov}}(\hat{\beta}_2, \hat{\beta}_1) & \widehat{\text{Var}}(\hat{\beta}_2) & \widehat{\text{Cov}}(\hat{\beta}_2, \hat{\beta}_3) \\ \widehat{\text{Cov}}(\hat{\beta}_3, \hat{\beta}_1) & \widehat{\text{Cov}}(\hat{\beta}_3, \hat{\beta}_2) & \widehat{\text{Var}}(\hat{\beta}_3) \end{bmatrix} \\ &= \begin{bmatrix} 0.3970573440 & -0.0066861059 & 0.0011014411 \\ -0.0066861059 & 0.1618026250 & -0.0038059631 \\ 0.0011014411 & -0.0038059631 & 0.1158848865 \end{bmatrix}.\end{aligned}$$

The following **R** script uses the **mvtnorm** package to simulate the multivariate normal data:

```
# include library to simulate correlated MV norm
library("mvtnorm")
```

* An alternate approach would be to simulate correlated random values drawn from a beta distribution which is constrained on the interval $[0, 1]$, with shape parameters α and β determined by estimated parameters and variances of those estimates. Computationally, this can be done by first generating standard normal variates with the required covariance structure, and then transforming them to beta variates with the required mean and standard deviation. See dos Santos Dias *et al.* (2008).

Now, we set up the parameter values needed to specify the simulation:

```
# number of samples to take from MV norm
iter <- 1000000;

# dipper parameter values and var-covar -- on the logit scale -- to use in simulation
beta1 <- 0.4512441; beta2 <- -0.1673372; beta3 <- -0.0159047;

vc <- matrix(c(0.3970573440,-0.0066861059,0.0011014411,
              -0.0066861059,0.1618026250,-0.0038059631,
              0.0011014411,-0.0038059631,0.1158848865),3,3,byrow=T);

# generate rannor samples conditional betas and VC matrix
logit_samples = rmvnorm(iter,mean=c(beta1,beta2,beta3),sigma=vc,method="svd")
```

Now, we check to confirm our simulated variance-covariance matrix is close to the estimated matrix (above):

```
# check to confirm simulated VC is correct...
cat("simulated VC matrix")
print(round(cov(logit_samples),10))

simulated VC matrix
      [,1]      [,2]      [,3]
[1,] 0.396846751 -0.006892055 0.001056857
[2,] -0.006892055 0.161725609 -0.003646518
[3,] 0.001056857 -0.003646518 0.115931880
```

Then, we simply back-transform our samples from the logit \rightarrow real probability scale, and proceed from there.

```
# convert logit samples to data frame
logit_samples <- as.data.frame(logit_samples)

# back-transform from logit scale
real_samples <- exp(logit_samples)/(1+exp(logit_samples));

# generate the product of back-transformed deviates
real_samples$prod = real_samples[,1]*real_samples[,2]*real_samples[,3];

# summary stats
cat("expectation of product =", mean(real_samples$prod))
cat("variance of product =", var(real_samples$prod))
```

Running this script results in the following estimates, which are quite close to the expected product (0.138871), and variance of the product derived using the Delta method (0.0025565):

```
expectation of product = 0.1370664
variance of product = 0.002377359
```

iii. using MCMC...

Another approach makes use of the Markov Chain Monte Carlo (MCMC) capabilities in **MARK**. Here, we provide only a brief description of the idea, and mechanics – for a more complete discussion, see Appendix E.

The basic idea is as follows. We'll fit model $\{\varphi_t p.\}$ to the male Dipper data, and use MCMC to derive estimates of the survival and encounter parameters, based on estimated moments (mean, median, or mode), and associated variances, from the posterior distribution for each of the parameters. The posterior distribution for each parameter is generated by Markov sampling over the joint probability distribution for all parameters, given the data.

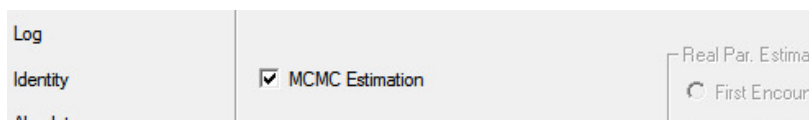
If we were using a specialized MCMC application, like **JAGS**, or **BUGS**, we could simply create a derived parameter as a function of other structural parameters in the model (say, $\mathbf{prod} = \varphi_1 \times \varphi_2 \times \varphi_3$), and then analyze the posterior samples for this derived parameter (this ability to explicitly code functions of parameters is one of the real conveniences of using MCMC, typically in a Bayesian framework).

The MCMC capabilities in **MARK** do not allow the explicit construction of a user-specified derived parameter. However, we can accomplish much the same thing, albeit in a slightly more 'brute-force' way, by simply

- (i) taking the individual sample chains from the MCMC simulations for each of the 3 real parameters involved in the 'function' of interest (i.e., the product – $\varphi_1 \rightarrow \varphi_3$),
- (ii) deriving the function of these parameters over the chains – in this case taking their product, and finally,
- (iii) evaluating this product as the posterior distribution for the product (which it is). In fact, this is equivalent to what **JAGS** or **BUGS** does, except that instead of calculating the product of the survival parameters at each step of the sampler, we simply do it *post hoc* – after the samplers are finished.

OK, let's see how this is done. First, we fit model $\{\varphi_t p.\}$ to the male Dipper data. We'll use a logit link (for reasons discussed in Appendix E).

Next, re-run this model. But, before submitting the model for numerical estimation for the second time, we first check the '**MCMC Estimation**' box:



Once you click the '**OK to Run**' button, **MARK** will respond with a window (shown at the top of the next page) where you specify the MCMC parameters that will specify aspects of the numerical estimation (see Appendix E for a complete discussion of these parameters).

What is generally important is that we want a sufficient number of samples (at all stages) to ensure that the samplers have converged on the stationary joint distribution. For this example we've used 7,000 'tuning' samples, 3,000 'burn in' samples, and 100,000 samples from the posterior distribution. We've also specified only a single chain, with no convergence diagnostics.

Once finished, **MARK** will output the results to the editor. If you scroll down to near the bottom of the output listing, you'll see various macro values that can be used for post-processing of the chains for each parameter. These macro values are copied into **SAS** or **R** programs that are provided in the **MARK** helpfile. We'll demonstrate the mechanics using **R**.

For the male Dipper data, and model $\{\varphi_t p.\}$, the **R** macro values are:

```
ncovs <- 7; # Number of beta estimates
nmeans <- 0; # Number of mean estimates
ndesigns <- 0; # Number of design matrix estimates
nsigmas <- 0; # Number of sigma estimates
nrhos <- 0; # Number of rho estimates
nlogit <- 7; # Number of real estimates
filename <- "C:\\USERS\\USER\\DESKTOP\\MCMC.BIN"; # path MCMC.BIN file
```

So, all we do is copy this into the appropriate section at the top of the **R** script provided in the **MARK** helpfile. The script is fairly lengthy, so we won't reproduce it in full here. Instead we'll focus on the additional steps you'll need to execute in order to derive an estimate of the variance for the product of the first 3 survival estimates.

First, copy the macro variables (above) into the **R** script, and execute it 'as is'. This will create an MCMC 'object', called '**mcmcdata**', that is compatible with one of several **R** packages (e.g., **coda**). This object contains each of the individual Markov chains, for each parameter.

Normally, what you'd do at this point is use some package, like **coda**, to post-process the chains, and generate various descriptive statistics and associated graphics. However, what we want to do here is estimate the variance of the product of $(\varphi_1 \times \varphi_2 \times \varphi_3)$. As outlined earlier, we will (i) extract the chains for the survival parameters φ_1 , φ_2 and φ_3 from the **mcmcdata** object, (ii) take their product, and (iii) generate various descriptive statistics for this product.

While there are a number of ways you might do this in **R**, the following works well enough. The first thing we do is convert the MCMC 'object' (**mcmcdata**) to a dataframe, which we'll call '**chaindata**':

```
chaindata <- as.data.frame(mcmcdata);
```

Next, we'll add a column to this new dataframe for the product ($\varphi_1 \times \varphi_2 \times \varphi_3$), and label this new column '**prod**'. Note, in the dataframe, these parameters are referred to (by their column names, which are explicitly set by the preceding R script) as '**real1**', '**real2**', and '**real3**', respectively:

```
chaindata$prod <- chaindata$real1*chaindata$real2*chaindata$real3;
```

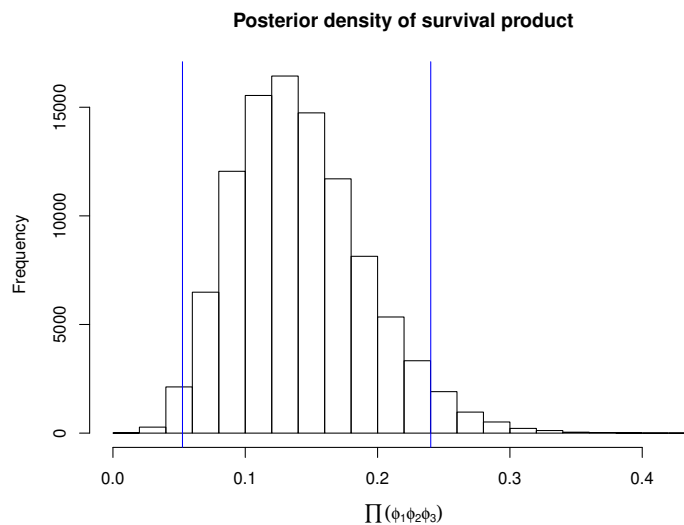
All that's left is to look at the summary statistics for this product:

```
summary(chaindata$prod)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
0.02155 0.10520 0.13620 0.14140 0.17180 0.47390

var(chaindata$prod)
[1] 0.002485553
```

As with the preceding approach, this script on the MCMC data results in moment estimates which are quite close to the expected product (0.138871), and variance of the product derived using the Delta method (0.0025565).

Couple of things to note. First, because the posterior for '**prod**' is asymmetrically distributed around the most frequent value (as shown in the following histogram), the median of this distribution, 0.13620, is perhaps the most appropriate moment to characterize the posterior, and is quite close to the expected product – somewhat closer than the mean.



Second, from a Bayesian perspective, it might be more meaningful to consider the credible interval, rather than the point estimate. Because of the asymmetry of the posterior distribution, use of the HPD (highest posterior density – see Appendix E) might be the most appropriate way to specify the interval. The upper and lower bounds of the 95% HPD for '**prod**', [0.05256, 0.24017], are plotted as vertical blue lines in the preceding frequency histogram.

Addendum B.2 – approximating $E[(X - \mu)^3]$ and beyond...

Here, we demonstrate the mechanics for evaluating $E((X - \mu)^3)$ (introduced in the -sidebar- on p. B - 15). While it might look daunting, in fact it is relatively straightforward, and is a convenient demonstration of how you can use the ‘algebra of moments’ to derive some interesting and useful things:

$$\begin{aligned}
 E((X - \mu)^3) &= \int_a^b (x - \mu)^3 p(x) dx \\
 &= \int_a^b (x^3 + 3\mu^2 x - 3\mu x^2 - \mu^3) p(x) dx \\
 &= \int_a^b x^3 p(x) dx + \int_a^b 3\mu^2 x p(x) dx - \int_a^b 3\mu x^2 p(x) dx - \int_a^b \mu^3 p(x) dx \\
 &= \int_a^b x^3 p(x) dx + 3\mu^2 \int_a^b x p(x) dx - 3\mu \int_a^b x^2 p(x) dx - \mu^3 \int_a^b p(x) dx \\
 &= M_3 + 3\mu^2(M_1) - 3\mu(M_2) - \mu^3(M_0).
 \end{aligned}$$

Since $M_1 = \mu$, and $M_0 = 1$, then

$$\begin{aligned}
 E((X - \mu)^3) &= \int_a^b (x - \mu)^3 p(x) dx \\
 &= M_3 + 3\mu^2(M_1) - 3\mu(M_2) - \mu^3(M_0) \\
 &= M_3 + 3M_1^3 - 3M_1M_2 - M_1^3.
 \end{aligned}$$

At this point, all that remains is substituting in the expressions for the moments corresponding to the particular pdf (in this case, $\mathcal{U}(a, b)$, as derived earlier in this Appendix), and you have your function for the expectation $E((X - \mu)^3)$.

We’ll leave it to you to confirm the algebra – the ‘answer’ is

$$\begin{aligned}
 E((X - \mu)^3) &= 2 \left(\frac{1}{2}a + \frac{1}{2}b \right)^3 - 3 \left(\frac{1}{2}a + \frac{1}{2}b \right) \left(\frac{1}{3}a^2 + \frac{1}{3}ab + \frac{1}{3}b^2 \right) + \frac{1}{4}a^3 + \frac{1}{4}a^2b + \frac{1}{4}ab^2 + \frac{1}{4}b^3 \\
 &= \text{‘some algebraic rearrangements, simplifications, and term cancellations’} \\
 &= 0.
 \end{aligned}$$

Yup, that’s it. A fair bit of work for what seems like an entirely anti-climatic result:

$$“E((X - \mu)^3) \text{ for the pdf } \mathcal{U}(a, b) \text{ is } 0.”$$

But you’re happier knowing how it’s done (no, really). We use the same procedure for $E((X - \mu)^4)$, and so on.

In fact, if you go through the exercise of calculating $E((X - \mu)^n)$ for $n = 4, 5, \dots$, you’ll find that they generally alternate between 0 (when n is odd), and non-zero (when n is even). This is a general result for any *symmetric* distribution (uniform, normal,...), and can be very helpful in simplifying the Taylor expansion. For *asymmetric* distributions (log-normal, exponential, gamma,...) you’ll need to include all the expectation terms.