

Count signal standard errors

Alex Reinhart and Nat DeFries

November 2021

Contents

1	Notation	1
2	Current status	2
2.1	The current code	2
2.2	The documented calculation	2
2.3	Derivation from covid-19 repo's aggregations-setup.R	3
3	Deriving the right approach	4
3.1	Option 1: Proceed with stated weight	5
3.2	Option 2: Proceed with different weight	5
3.3	Option 3: Forget a term	6
4	Conclusion	7

1 Notation

Let n be the number of observations. Let:

$Y_i \in [0, 1]$ = respondent i 's response

w_i = respondent i 's weight.

We normalize the weights so that $\sum_{i=1}^n w_i = 1$. Note that Y_i is a proportion, since each respondent reports the proportion of their household that is sick.

2 Current status

2.1 The current code

In the current code (`count.R`), the function `compute_count_response` calculates the following:

$$\hat{p} = \sum_{i=1}^n w_i Y_i \quad (1)$$

$$\hat{s}e_{\text{unadj}} = \sqrt{\sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2} \quad (2)$$

$$n_{\text{eff}} = n \frac{\left(\frac{1}{n} \sum_{i=1}^n w_i\right)^2}{\frac{1}{n} \sum_{i=1}^n w_i^2} = \frac{(\sum_{i=1}^n w_i)^2}{\sum_{i=1}^n w_i^2} = \frac{1}{\sum_{i=1}^n w_i^2}. \quad (3)$$

(Evidently the n_{eff} code can be simplified in the knowledge that the weights sum to 1.)

Then `jeffreys_se` adjusts this to produce a final weighted standard error by introducing an additional observation. This additional observation has $Y = 0.5$. It does this:

$$\hat{s}e_{\text{code}} = \frac{1}{1 + n_{\text{eff}}} \sqrt{\left(\frac{1}{2} - \hat{p}\right)^2 + n_{\text{eff}}^2 \hat{s}e_{\text{unadj}}^2}. \quad (4)$$

2.2 The documented calculation

The documentation produces the same \hat{p} as the code. It then shows that:

$$\hat{s}e_{\text{unadj}} = \sqrt{\sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2}, \quad (5)$$

which is the same as in the code as well. But it then states the following adjusted standard error:

$$\hat{s}e_{\text{doc}} = \sqrt{\left(\frac{1}{1 + n_{\text{eff}}}\right)^2 \left(\frac{1}{2} - \hat{p}\right)^2 + n_{\text{eff}} \hat{s}e_{\text{unadj}}^2}, \quad (6)$$

again using the same definition of n_{eff} .

So how is $\hat{s}e_{\text{code}}$ different from $\hat{s}e_{\text{doc}}$? Let's try to pull factors out of $\hat{s}e_{\text{doc}}$ to make the comparison clearer:

$$\hat{s}e_{\text{doc}} = \frac{1}{1 + n_{\text{eff}}} \sqrt{\left(\frac{1}{2} - \hat{p}\right)^2 + (1 + n_{\text{eff}})^2 n_{\text{eff}} \hat{s}e_{\text{unadj}}^2}. \quad (7)$$

That $(1 + n_{\text{eff}})^2$ term is what Nat found to be inconsistent with the code, since expanding it out leads to several terms not present in the code. (Nat's derivation in the GitHub comments started from (6) but without the n_{eff} multiplying $\hat{s}e_{\text{unadj}}^2$, since he thought that to be a typo.)

2.3 Derivation from covid-19 repo's aggregations-setup.R

The `jeffreys_se` function in `aggregations-setup.R` contains a derivation in code comments, which we translate into notation here. For the purposes of this derivation, let's **not** assume that the weights are normalized. We'll apply assumptions at the end.

In the general case,

$$\hat{\text{se}}_{\text{unadj}} = \sqrt{\frac{\sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2}{(\sum_{j=1}^n w_j)^2}} \quad (8)$$

$$n_{\text{eff}} = \frac{(\sum_{i=1}^n w_i)^2}{\sum_{i=1}^n w_i^2}. \quad (9)$$

Rearranging,

$$\left(\sum_{j=1}^n w_j\right)^2 \hat{\text{se}}_{\text{unadj}}^2 = \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2. \quad (10)$$

Applying the Jeffreys adjustment on both sides by adding an extra term to the sums with $Y_0 = 0.5$ and weight w_0 ,

$$\left(w_0 + \sum_{j=1}^n w_j\right)^2 \hat{\text{se}}_{\text{final}}^2 = w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2, \quad (11)$$

and rearranging, we obtain

$$\hat{\text{se}}_{\text{final}} = \sqrt{\frac{w_0^2}{(w_0 + \sum_{j=1}^n w_j)^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \frac{1}{(w_0 + \sum_{j=1}^n w_j)^2} \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2} \quad (12)$$

$$= \sqrt{\frac{w_0^2}{(w_0 + \sum_{j=1}^n w_j)^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \frac{(\sum_{j=1}^n w_j)^2}{(w_0 + \sum_{j=1}^n w_j)^2} \hat{\text{se}}_{\text{unadj}}^2} \quad (13)$$

$$= \frac{1}{w_0 + \sum_{j=1}^n w_j} \sqrt{w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \left(\sum_{j=1}^n w_j\right)^2 \hat{\text{se}}_{\text{unadj}}^2}. \quad (14)$$

Assuming the weights are normalized such that $\sum_{j=1}^n w_j = 1$, this becomes

$$\hat{\text{se}}_{\text{final, normalized}} = \frac{1}{w_0 + 1} \sqrt{w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{\text{se}}_{\text{unadj}}^2} \quad (15)$$

Now we must choose the weight w_0 for the pseudo-observation. Let $w_0 = 1/n_{\text{eff}}$. This

results in the following:

$$\hat{s}e_{\text{final,normalized}} = \frac{1}{1/n_{\text{eff}} + 1} \sqrt{\frac{1}{n_{\text{eff}}^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{s}e_{\text{unadj}}} \quad (16)$$

$$= \frac{n_{\text{eff}}}{1 + n_{\text{eff}}} \sqrt{\frac{1}{n_{\text{eff}}^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{s}e_{\text{unadj}}} \quad (17)$$

$$= \frac{1}{1 + n_{\text{eff}}} \sqrt{\left(\frac{1}{2} - \hat{p}\right)^2 + n_{\text{eff}}^2 \hat{s}e_{\text{unadj}}^2}. \quad (18)$$

This matches the current implementation in (4).

3 Deriving the right approach

So what's a typo and what's not a typo? Let's start over.

We can agree that

$$\hat{s}e_{\text{unadj}} = \sqrt{\frac{\sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2}{(\sum_{j=1}^n w_j)^2}}, \quad (19)$$

and that we want to prevent this from being zero, which would occur when $Y_i = \hat{p}$ for all i such that $w_i > 0$. (For example, if $\hat{p} = 0$.)

There are two ways to achieve this that would fit the Jeffreys idea:

1. Transform \hat{p} : recalculate it as if there were one additional observation $i = 0$ with $Y_0 = 1/2$, then calculate the standard error according to the equation above while including that extra observation in the sum.
2. Leave \hat{p} untouched but add the extra term to the sum anyway.

We chose the latter route because we expected $\hat{p} \approx 0$ for CLI signals, and Jeffreys adjustment of \hat{p} would cause large amounts of bias for values very close to zero.

So we want to add an extra term to the sum. But what should its weight w_0 be? Here the documentation is not specific: it says this extra observation has “weight assigned to appear like a single effective observation according to importance sampling diagnostics.”

Let's just use w_0 and see what values of w_0 match the code or documentation. If we

start with (19), we have

$$\left(\sum_{i=1}^n w_i\right)^2 \hat{\text{se}}_{\text{unadj}}^2 = \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2 \quad (20)$$

$$\left(w_0 + \sum_{i=1}^n w_i\right)^2 \hat{\text{se}}_{\text{final}}^2 = w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2 \quad (21)$$

$$\hat{\text{se}}_{\text{final}} = \frac{1}{w_0 + \sum_{i=1}^n w_i} \sqrt{w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2} \quad (22)$$

3.1 Option 1: Proceed with stated weight

If

$$w_0 = \frac{1}{1 + n_{\text{eff}}}, \quad (23)$$

and $\sum_{i=1}^n w_i = 1$, then starting with (22), we obtain:

$$\hat{\text{se}}_{\text{final}} = \frac{1}{\frac{1}{1+n_{\text{eff}}} + 1} \sqrt{\left(\frac{1}{1+n_{\text{eff}}}\right)^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2} \quad (24)$$

$$= \frac{1}{\frac{1}{1+n_{\text{eff}}} + 1} \sqrt{\left(\frac{1}{1+n_{\text{eff}}}\right)^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{\text{se}}_{\text{unadj}}^2} \quad (25)$$

$$= \frac{1}{\frac{1+1+n_{\text{eff}}}{1+n_{\text{eff}}}} \sqrt{\left(\frac{1}{1+n_{\text{eff}}}\right)^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{\text{se}}_{\text{unadj}}^2} \quad (26)$$

$$= \frac{1+n_{\text{eff}}}{2+n_{\text{eff}}} \sqrt{\left(\frac{1}{1+n_{\text{eff}}}\right)^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{\text{se}}_{\text{unadj}}^2} \quad (27)$$

$$= \frac{1}{2+n_{\text{eff}}} \sqrt{\left(\frac{1}{2} - \hat{p}\right)^2 + (1+n_{\text{eff}})^2 \hat{\text{se}}_{\text{unadj}}^2}. \quad (28)$$

This doesn't match either documentation or code.

3.2 Option 2: Proceed with different weight

If

$$w_0 = \frac{1}{n_{\text{eff}}}, \quad (29)$$

and $\sum_{i=1}^n w_i = 1$, then starting with (22), we obtain:

$$\begin{aligned}\hat{s}e_{\text{final}} &= \frac{1}{1/n_{\text{eff}} + 1} \sqrt{\frac{1}{n_{\text{eff}}^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2} \\ &= \frac{n_{\text{eff}}}{1 + n_{\text{eff}}} \sqrt{\frac{1}{n_{\text{eff}}^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2} \\ &= \frac{1}{1 + n_{\text{eff}}} \sqrt{\left(\frac{1}{2} - \hat{p}\right)^2 + n_{\text{eff}}^2 \hat{s}e_{\text{unadj}}^2},\end{aligned}$$

which matches $\hat{s}e_{\text{code}}$ in (4).

3.3 Option 3: Forget a term

We started from $\hat{s}e_{\text{unadj}}$ in (19). But suppose we started by normalizing the weights to sum to 1, meaning we started with this version instead:

$$\hat{s}e_{\text{unadj}} = \sqrt{\sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2}. \quad (30)$$

Notice the absence of a denominator. When we proceed to make the Jeffreys version, that denominator does not get moved to the left-hand side, and so when we add the fake term with $Y_0 = 0.5$ and weight w_0 , that term is **not** added to the left-hand side. Instead, we get this:

$$\hat{s}e_{\text{final}}^2 = w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \sum_{i=1}^n w_i^2 (Y_i - \hat{p})^2 \quad (31)$$

$$\hat{s}e_{\text{final}} = \sqrt{w_0^2 \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{s}e_{\text{unadj}}^2} \quad (32)$$

If we now substitute in $w_0 = 1/(1 + n_{\text{eff}})$, we get

$$\hat{s}e_{\text{final}} = \sqrt{\frac{1}{(1 + n_{\text{eff}})^2} \left(\frac{1}{2} - \hat{p}\right)^2 + \hat{s}e_{\text{unadj}}^2} \quad (33)$$

$$= \frac{1}{1 + n_{\text{eff}}} \sqrt{\left(\frac{1}{2} - \hat{p}\right)^2 + (1 + n_{\text{eff}})^2 \hat{s}e_{\text{unadj}}^2}. \quad (34)$$

then this matches the documented $\hat{s}e_{\text{doc}}$ —except for the documentation’s excess n_{eff} term that Nat identified as a typo.

4 Conclusion

The documentation contains two mistakes: it has an extra n_{eff} , and it did not include the denominator in $\hat{\text{se}}_{\text{unadj}}$, causing it to miss an extra term.

The code does not match the documentation's claim that $w_0 = 1/(1 + n_{\text{eff}})$. It avoids the missing denominator, but defines $w_0 = 1/n_{\text{eff}}$.

The question is then “What should w_0 be?” To resolve this, let's consider the case where all weights are equal, i.e. $w_i = w$. Then

$$n_{\text{eff}} = \frac{n^2 w^2}{n w^2} = n. \quad (35)$$

If all weights are equal and sum to 1 before the fake observation is added, $w = 1/n$. it only makes sense for $w_0 = w$. We could only obtain this result if $w_0 = 1/n_{\text{eff}} = 1/n$.

If all weights are equal and sum to 1 *after* the fake observation is added, $w = 1/(n+1)$, i.e. $w_0 = 1/(n_{\text{eff}} + 1)$.

The code normalizes before the fake observation is added, so it makes sense that $w = 1/n_{\text{eff}}$.

In conclusion, the code is correct and the documentation is wrong. Option 2 is the way to go.