

## Review of Wang and Oskin, 'Combined linear regression...'

---

**Summary.** As noted by some of the reviewers, to some extent this paper is a solution in search of a problem. Cosmogenic-nuclide depth profiles are usually interpreted by inversion of a forward model that predicts nuclide concentrations and whose parameters are the age of a surface, the erosion rate of a surface, and some nuisance parameters typically including inherited nuclide concentrations. Generally this approach works fine, or as well as can be expected given the inherent lack of age resolution for typical depth profiles in which the muon-produced inventory is small in relation to the inheritance.

This paper presents a nice demonstration that you can perform a very simple inversion of depth-profile data for age and inheritance using linear regression, if you change coordinates from depth to production rate. This does help to simplify the problem somewhat and make it more easily accessible, although probably not as much as suggested by the authors. Thus, from this perspective I think this is a contribution that is certainly of interest for publication. However, there are some issues that have come up in review that I do think need further attention, as described below.

**Why needed?** One of the issues that has come up in the discussion of this paper is that from the applications perspective there is not a strong need for a simplified inversion. Depth-profile data are not collected in great quantity and there is not a science application that is currently seriously hindered by the computational time needed to do a full forward model inversion. For myself I am not worried about this issue and I don't think it's an obstacle to publication. For one thing, it is potentially useful for making sure that a more complex inversion scheme is working correctly. Also, I can envision a fast inversion method being useful for database applications in which one seeks to compare a lot of depth-profile results using different production rate scaling methods, or something of that nature. Of course it's actually not that fast because you still need to estimate all the production rates due to muons, which requires a site-specific muon production calculation, which in turn requires a bunch of numerical integrations no matter what. Regardless, however, even though a simplified inversion method isn't of dire immediate need for any current application, it is certainly something that is potentially useful.

**The issue of negative inheritance.** On the other hand, a second issue that has been the focus of much of the review discussion seems to be a more serious problem. This has to do with error propagation and generating an uncertainty distribution for the surface exposure age. It is not really feasible to come up with an analytical expression for the uncertainty in  $T_e$  and  $C_{inh}$ , or even to use something like a York regression, because many of the input uncertainties (e.g., in  $P_{zi}$ , or in the mass depths because they all depend on the same density measurements) are correlated in a complicated way. Thus, the authors use a Monte Carlo simulation where they vary the input parameters and carry out the regression many times to generate an uncertainty distribution. This brings up the issue, which has been discussed at length in the review comments, of how to handle the fact that many Monte Carlo realizations generate negative values for nuclide inheritance. The authors have proposed, and discuss in both paper drafts and the response to reviews, two methods for dealing with this: first, accepting all Monte Carlo results even if they yield unphysical negative values of the inheritance; second, discarding as unphysical Monte Carlo realizations that yield negative values for the inheritance.

Here I will argue that both of these approaches are incorrect. Although this overall subject seems like a rather arcane point, similar situations often occur in cosmogenic-nuclide applications where a forward model that includes inheritance is being fit to data. Thus, I think airing this issue in the discussion of this paper is valuable and helpful for the field overall.

I think the reason that both approaches are incorrect is that a simple linear regression for  $C_{inh}$  and  $T_e$  given  $C = PT_e + C_{inh}$  (here I am just abbreviating Equation 4a by replacing all the production terms with  $P$ ) is an incomplete description of the regression problem. In other words, even if one leaves aside the question of whether linear regression is the best way to determine the age, the authors have not defined the regression problem itself appropriately.

Basically, linear regression is a least-squares optimization problem:

Problem 1:

$$\text{given} \quad X = (T_e, C_{inh}) \quad (1a)$$

$$\text{minimize} \quad f(X) = \sum_i [(C_{inh} + P_i T_e) - C_i]^2 \quad (1b)$$

$$\text{over } X \text{ such that: } -\infty < T_e < \infty; -\infty < C_{inh} < \infty \quad (1c)$$

$$(1d)$$

This is how the authors are treating it, as an unconstrained optimization problem in which you can optimize over all  $X$ , meaning any possible values of  $T_e$  and  $C_{inh}$ . This leads to the normal simple formulae for linear regression.

Unfortunately this is an incomplete description of the problem. Because of the physical constraints that age and inheritance can't be negative, the problem is actually something different:

Problem 2:

$$\text{given} \quad X = (T_e, C_{inh}) \quad (2a)$$

$$\text{minimize} \quad f(X) = \sum_i [(C_{inh} + P_i T_e) - C_i]^2 \quad (2b)$$

$$\text{over } X \text{ such that: } 0 < T_e < \infty; 0 < C_{inh} < \infty \quad (2c)$$

$$(2d)$$

Problem 1 and Problem 2 are equivalent only if the constraints are inoperative, which occurs when  $C_{inh} \gg 0$  and  $T_e \gg 0$ . If that's not the case, then the problems are not equivalent. Solving the unconstrained problem won't necessarily also give the answer to the constrained problem.

Unfortunately the constrained problem requires a numerical solution. Non-negative least squares problems are common enough to have a Wikipedia entry, so there are canned algorithms in many programming environments. Regardless, this causes some trouble in the context of the paper because a key point of the paper is that linear regression can be done with a simple formula, not a numerical optimization. That is, you can solve Problem 1 just by writing down the linear regression formula. If you add constraints to the optimization, this isn't true any more – you have to use a numerical optimization scheme. So if you have to consider the complete problem – Problem 2 – you have to use a numerical method anyway and much of the simplicity advantage of linear regression disappears.

Of course, the difference between Problem 1 and Problem 2 mostly doesn't affect inverting the depth profile once by linear regression – the initial attempt would only fail if the initial regression from the data as measured gives  $T_e < 0$  or  $C_{inh} < 0$ , which isn't going to happen in most applications. It only affects figuring out what the uncertainty distribution is in a Monte Carlo simulation. If either the age or the uncertainty are near zero, then solving the unconstrained problem repeatedly and discarding all iterations that yield  $C_{inh} < 0$  (or  $T_e < 0$ , which is not mentioned in the paper but could happen) will NOT give the same result as solving the constrained problem repeatedly.

Consider three possibilities:

- Option 1. Solve Problem 1 in each iteration.
- Option 2. Solve problem 1 in each iteration, but discard all iterations in which  $C_{inh} < 0$ .
- Option 3. Solve problem 2, the constrained optimization, using a numerical optimizer in each Monte Carlo iteration.

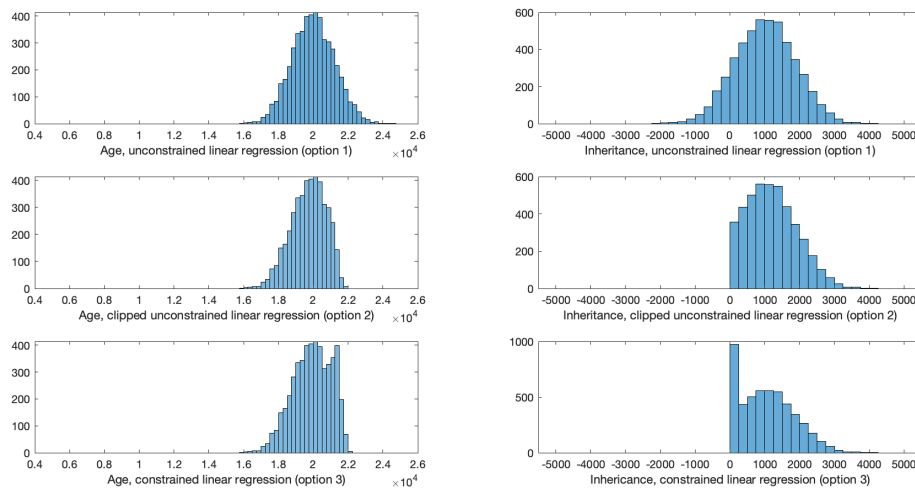


Figure 1: *Monte Carlo uncertainty estimates for a simplified regression problem with an age of 20 ka and inheritance equal to 1 ka of exposure.*

If the constraints are inoperative ( $C_{inh} \gg 0$  and  $T_e \gg 0$ ) these all give the same result. However, consider Figure 1, which shows Monte Carlo inversion results for a simple test problem where the age is 20,000 years and the inheritance is equal to 1000 years exposure, so the constraints are operative. Option 1 (unconstrained, no clipping) results in normal uncertainty distributions for concentration and age, but the uncertainty distribution for inheritance incorrectly extends to negative values. Option 2 removes all iterations that yield negative inheritance, which, obviously, yields a clipped normal distribution for inheritance, but less obviously produces a low-skewed distribution for the age, because the steeper slopes that led to negative intercepts on the inheritance axis are discarded. Option 3, the correct constrained optimization approach, forces all values of the inheritance that would have been negative in the unconstrained problem to be exactly zero, which leads to a secondary mode at the high end of the age distribution. Thus, both Option 1 and Option

2 lead to incorrect estimates of the uncertainty distribution. In this example, Option 1 incorrectly indicates that there is a finite likelihood that the age is greater than  $\sim 22$  ka. Option 2 incorrectly underestimates the likelihood that the age is between  $\sim 20$  and  $\sim 22$  ka.

This is, in fact, rather weird behaviour, which maybe is a hint that linear regression is maybe not the best way to generate uncertainty distributions when age or inheritance is close to zero. It is certainly possible that what we should be learning here is that if we want an uncertainty distribution for the age when inheritance is close to zero, we should use something different, possibly more like a Bayesian fitting scheme with a prior restricted to positive age and inheritance.

So, the authors are correct in remarking in their response that only Option 1 leads to a normal uncertainty distribution. However, their assertions that “negative results, though physically impossible, are necessary for mathematical reasons” and “the overall distributions of the inheritance should be centered around 0, meaning that approximately half of the estimated inheritance should be negative” are, in my opinion, not correct. In fact, the physical requirement is that NONE of the estimated inheritance values (or the age) should be negative, which is why the regression problem is required to be a constrained and not an unconstrained optimization. This is not expected to lead to a normal or symmetrical uncertainty distribution. It implies non-normal uncertainty distributions for both the age and inheritance (which, by the way, should not be represented by means and standard deviations, but rather statistics that aren’t specific to a particular distribution, for example mode and confidence intervals).

This effect is more striking (although less pathological-looking) if the inheritance really is close to zero. Figure 2 shows the same results for an age of 10 ka and zero inheritance. Here the correct uncertainty distribution for the age (Option 3) is extremely skewed and of course representing it as a mean and standard deviation would be very misleading.

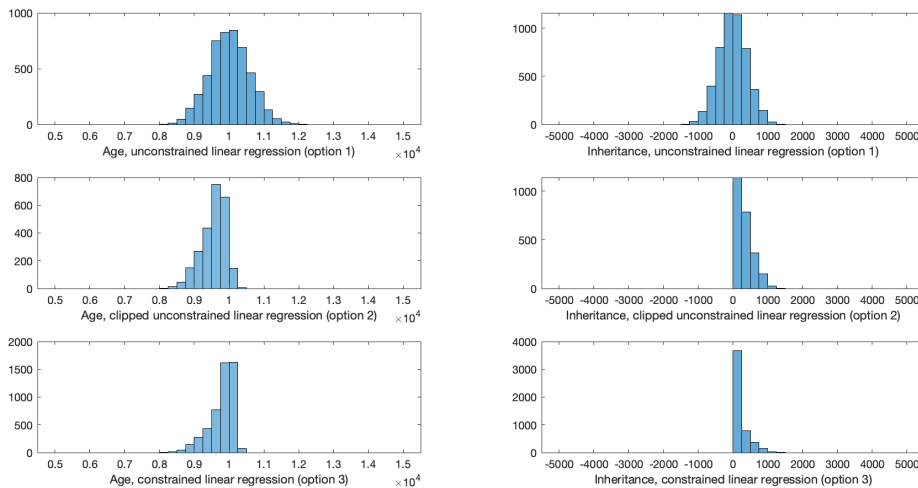


Figure 2: Monte Carlo uncertainty estimates for a simplified regression problem with an age of 10 ka and zero inheritance.

To summarize, the question here isn’t really which uncertainty distribution is correct (I would say that

once you have chosen linear regression as the overall method, only option 3 is correct), but whether the uncertainty distributions generated by either Option 1 or Option 2 are close enough to the correct answer that the inaccuracy incurred by solving an oversimplified problem is outweighed by the convenience of using the simple unconstrained regression formula. Basically, the answer to this is that if  $C_{inh} \gg 0$  and  $T_e \gg 0$ , the uncertainty distribution derived from the unconstrained problem is totally fine. If this isn't true, the unconstrained problem gives the wrong uncertainty distribution, although it might not be that different, so this might not matter in many applications.

In summary, I apologize for discussing at great length the issue of negative inheritance, which is, in fact, an arcane side issue that only affects the uncertainty analysis and is not the main point of this paper. However, as noted, similar issues come up in other cosmogenic-nuclide applications and I think a thorough discussion of the issue is helpful. Regardless, I think the regression model discussed in this paper is interesting and potentially useful, and I am supportive of publication. However, I think the authors do need to revise discussion of the negative inheritance issue. Specifically,

1. As written, several parts of section 3.1.2 are not correct, for example “imposing the physically reasonable prerequisite...may lead to underestimation of the exposure age.” It's not a physically reasonable but optional prerequisite, it's a requirement. It's also inaccurate to say that it leads to underestimation of the exposure age – what it actually does is lead to an incorrect uncertainty distribution, and the lower value of the age appears to be the result of improperly using the mean to represent an asymmetrical distribution. In the constrained regression problem, the uncertainty distribution is not expected to be symmetrical about the true value. Thus, this section needs revision so that it correctly outlines how (i) the actual regression problem is a constrained linear regression that is expected to lead to complicated uncertainty distributions when the constraints are operative, and (ii) applying an unconstrained regression is a simplification that is only correct when  $C_{inh} \gg 0$  and  $T_e \gg 0$ .
2. It's also probably worth a brief discussion of the case where a surface is quite young, so that the uncertainty distribution for  $t$  runs into zero. The unconstrained regression is also inappropriate in this case. This could possibly occur even when  $t$  is fairly large if the inheritance is also large, such that  $T_e P < C_{inh}$ .
3. The Beida River / Fig 8 analysis, which improperly uses an unconstrained regression when constraints are operative, should be redone with the correct, constrained regression. Of course for the Lees Ferry analysis, the constraints are inoperative, so the unconstrained regression is fine.

In any case, items 1-3 above are the main changes that I think are needed for publication.

**Other items.** In addition, there are a few lesser items that should be corrected before publication.

The most important one is that the description of the effective attenuation lengths for muon production (line 62) is oversimplified and therefore somewhat misleading. There is no single attenuation length for either fast muon production or negative muon capture production, because the nature of the production process is such that as depth increases, the energy of the remaining muons increases, so the instantaneous attenuation length for the production process also increases. Thus, describing  $\Lambda_{m1}$  and  $\Lambda_{m2}$  as the attenuation lengths for these processes is not correct. The values of 1500 and 4320 g cm<sup>-2</sup> that are in Table 1 in this paper, which were given by Heisinger as approximate values that could be used in simplified erosion rate integrations, are not correct at any site or depth, except possibly by accident. Using these values in an application with significant muon production would most likely yield a result that was quite wrong. However, it is true (see Balco, 2019, section 8) that it is usually possible to represent total production by muons in a finite depth range at a specific site as the sum of two exponential functions. Although the authors' response to the reviews stated that “We

updated muon production rates used in the pseudo profiles and in the Beida River case,' this is not evident from the text, and as far as I can tell, it appears likely that what they actually did was modify the surface production rates and not the attenuation constants. More seriously, it is not possible to tell whether they used site-specific values for the Beida River and Lees Ferry example, or if they used the incorrect values in Table 1. At the very least, the authors should modify the text near line 62 to indicate that the values of  $P$  and  $\Lambda$  pertaining to muon production are not, in fact, true production rates or attenuation lengths for any particular process, but instead are site-specific constants that come from fitting to a more complicated production model.

Two lesser issues just pertain to reducing confusion in the paper:

Starting in section 3 of the paper the authors contrast the results of their regression scheme with a full forward-model-fitting scheme that they describe as 'Bayesian.' This is confusing because the important point is that this method employs a full forward model to predict the observables – the contrast between this method and the regression method would be the same whether the approach to choosing the best values of the model parameters approach was Bayesian, frequentist, or something else. Likewise, it would be possible to perform a Bayesian linear regression. Thus, it would be more helpful to the reader to describe this alternative method as 'forward model fitting' or 'forward model optimization' rather than 'Bayesian' by itself.

Near line 40, the authors use the term 'derivative' to describe refined approaches that are derived from the initial regression approach. The use of this word is confusing here, because in the context of a paper like this one that is mathematical in nature, it implies to the reader that these approaches will involve derivatives in the mathematical sense. For example, in line 44, the use of 'second derivative approach' indicates to the reader that the approach will involve second derivatives of some function or field, which is not the case. Thus, 'derivative' should not be used here. Possible improvements would be to use 'specific approaches' or 'special cases' to contrast with the 'general approach' in line 40.

Finally, a thoroughly insignificant point is that I disagree with reviewer Alan Hidy about the abbreviation 'TCN.' I see no reason that the method in the paper could not be applied to depth profiles in lunar or Martian regolith! The authors should use their best judgement here.