1    Response to Valle and Zorello Laporta: Clarifying the use of instrumental variable methods to

2    understand the effects of environmental change on infectious disease transmission

3

4    Running head: Response to the critique by Valle & Zorello Laporta

5

6    Andrew J. MacDonald[1*] & Erin A. Mordecai[2]

7    1. Earth Research Institute and Bren School of Environmental Science and Management,

8    University of California, Santa Barbara, CA, USA

9    2. Department of Biology, Stanford University, Stanford, CA, USA

10   *Corresponding author: Bren School of Environmental Science and Management, University of

11   California, Santa Barbara, CA 93106-5131; andy.j.macdon@gmail.com

12

15

16   Abstract: 164 words

17   Text: 2264 words

18   Figures: 0

19   Tables: 0

20   SI: 1 SI table

21

22

23    Abstract

24       Identifying the effects of environmental change on the transmission of vector-borne and

25    zoonotic diseases is of fundamental importance in the face of rapid global change. Causal

26    inference approaches, including instrumental variable (IV) estimation, hold promise in

27    disentangling plausibly causal relationships from observational data in these complex systems.

28    Valle and Zorello Laporta recently critiqued the application of such approaches in our recent

29    study of the effects of deforestation on malaria transmission in the Brazilian Amazon on the

30    grounds that key statistical assumptions were not met. Here, we respond to this critique by: 1)

31    deriving the IV estimator in order to clarify the assumptions that Valle and Zorello Laporta

32    conflate and misrepresent in their critique; 2) discussing these key assumptions as they relate to

33    our original study and how our original approach reasonably satisfies the assumptions; and 3)

34    presenting model results using alternative instrumental variables that can be argued more

35    strongly satisfy key assumptions, illustrating that our results and original conclusion—that

36    deforestation drives malaria transmission—remain unchanged.

37

38

39

40

41

42

43

44

45    Main Text

46          There is substantial and increasing interest in understanding the role that processes of

47    global change are playing in the ecology and transmission of vector-borne and zoonotic

48    diseases.[1,2] While these questions are of fundamental importance given the increasing rate of

49    climate and land use change, and the large proportion of emerging infectious diseases that are

50    vector-borne or of zoonotic origin,[3] causally linking these two processes is an enormous

51    challenge. Take as an example the case of deforestation impacts on malaria transmission in the

52    Brazilian Amazon, the focus of MacDonald & Mordecai[4] and the critique by Valle & Zorello

53    Laporta.[5] The gold standard of a randomized controlled trial in which deforestation is

54    experimentally manipulated and randomly assigned to different regions to assess its impact on

55    malaria transmission presents obvious logistical and ethical barriers that make such an approach

56    largely infeasible. As a result, researchers must rely on observational data and employ statistical

57    approaches to approximate, as closely as possible, the experimental ideal.

58          One promising set of statistical techniques—broadly referred to as causal inference

59    methods, which includes Instrumental Variable (IV) estimation, are increasingly being leveraged

60    to disentangle plausibly causal relationships from observational data in ecology. Due to the

61    challenges described above, these approaches have been employed by researchers assessing

62    global change impacts on infectious disease,[6-14] including in another recent study investigating

63    the effects of deforestation on malaria transmission in Brazil,[14] with similar results to our own

64    work. Valle and Zorello Laporta[5] rightly point out that model assumptions are critically

65    important in such approaches, and that causal conclusions should be carefully drawn in these

66    contexts. However, the authors unfortunately conflate the assumptions of IV estimation in their

67    perspective piece. As a relatively new approach in ecology and environmental science,[6] it is

68    important that the underlying assumptions are clear for appropriate application.

69        IV is a useful approach to overcome what is known as endogeneity bias, which is due to a

70    relationship between the error term and one or more of the explanatory variables, (formally,

71    $E[\varepsilon_i|x_i] \neq 0$ where $\varepsilon$ and $x$ represent the error term and explanatory variable for observation $i$).

72    Such a relationship could be due to bidirectional causality where, for example, deforestation may

73    drive malaria transmission but malaria burden may also influence rates of deforestation. In IV, a

74    third variable, known as an instrument ($z_i$), is used to isolate exogenous variation in explanatory

75    variable $x_i$ and recover a statistically consistent estimator for the true relationship between the

76    exogenous variable and the outcome.

77        The instrument must meet two conditions for IV to be a consistent estimator, which are

78    sometimes termed "relevance" and "exclusion" criteria. In words, the instrument must be

79    statistically associated with the endogenous variable ("relevance") and must be related to the

80    outcome only through its relationship with the endogenous variable ("exclusion"). While the

81    wording is easy to remember, it leaves much open to interpretation. For example, does relevance

82    require a causal link? Does exclusion require statistical independence? The derivation makes

83    these key assumptions much more apparent. Before showing the derivation, we will first provide

84    brief background to our original study,[4] the critique by Valle & Zorello Laporta[5] and our

85    response.

86        In MacDonald & Mordecai,[4] we were first interested in predicting annual malaria

87    incidence as a function of annual deforestation, and use aerosol optical depth (AOD) in the

88    month of September from MODIS satellite imagery as our "instrument." We expand on the

89    methodology and terminology below, but set the context of the argument here. Valle & Zorello

90    Laporta[5] have two critiques of our IV approach. The first, however, is a misrepresentation of the

91    assumptions of IV, namely that a valid IV requires that the IV has a *causal* effect on the

92    endogenous explanatory variable. They state, "However, it is deforestation that causes aerosol

93    pollution […] rather than aerosol pollution that causes deforestation […] As a result, [the

94    relevance] assumption is clearly violated." As we show below, causality is not required.[15]

95    Rather, there must be an "association", or more specifically, the covariance between the

96    instrument and the endogenous variable must not be zero. However, it is possible that an

97    instrumental variable itself introduces endogeneity bias if it does not meet the exclusion criteria,

98    and this can be particularly problematic in the case of "weak instruments" as we show below.

99    This can occur, for example, in cases where the instrument (e.g., AOD) is strongly driven by the

100    endogenous predictor variable (e.g., deforestation). In our case, we chose AOD as an instrument

101    for deforestation, as it is an indicator of human activity on the landscape.[16] Further, over our

102    study period, AOD was decoupled from deforestation as biomass burning in the Brazilian

103    Amazon—and resulting AOD—was primarily driven by fires intentionally set to keep *existing*

104    pastures and agricultural lands clear[16] and by drought conditions leading to wildfires in already

105    degraded forests,[16-18] rather than by new deforestation activity.

106         Nevertheless, to explore the extent to which our original IV estimates of the effect of

107    deforestation on malaria may have been affected by potential endogeneity introduced by the use

108    of AOD as an IV, we run additional IV models using 1) last year's AOD as an instrument for this

109    year's deforestation, and 2) remotely sensed, average municipality soil quality[19] processed in

110    Google Earth Engine,[20] interacted with annual international soy and beef commodity prices from

111    the World Bank. We chose last year's AOD since it is correlated with this year's deforestation

112    (relevance), but this year's deforestation could not have caused last year's AOD. While this

113    addresses the issue of reverse causality, it is plausible that there remain endogeneity issues in this

114    context. For example, if last year's AOD somehow acts upon this year's malaria through

115    mechanisms beyond deforestation, then the exclusion criteria would fail. To address these

116    potential lingering concerns, we run additional models using soil quality coupled with

117    international agricultural commodity prices for key Brazilian exports, which may influence a

118    land owners' decision to clear forest for agricultural production (relevance); in this case,

119    deforestation rates do not cause soil quality and are highly unlikely to shift international

120    commodity prices (exclusion). We run these IV models on our interior Amazon sample of

121    municipalities, where active deforestation rates are highest and where we predict forest clearing

122    should have the strongest effect on malaria transmission,[4] predicting both total malaria and

123    *Plasmodium falciparum* malaria incidence, following our original study.[4] Results are presented

124    in the SI (Table S1). In brief, we find significant positive effects of deforestation on malaria

125    transmission in each of these additional model specifications, with coefficients of similar, though

126    slightly larger magnitude than our original study. Our main conclusion, that deforestation

127    increases malaria transmission in the Brazilian Amazon, remains unchanged.

128            The second goal of MacDonald & Mordecai[4] is to understand whether annual malaria

129    burden feeds back to influence annual rates of deforestation, and we use optimal temperature for

130    malaria transmission in the dry season as our instrument for malaria. Optimal temperature was

131    defined as the sum of days falling within a narrow temperature band that is optimal for malaria

132    transmission (24-26ºC) based on earlier mosquito and parasite trait-based mechanistic modeling

133    studies.[21] Valle & Zorello Laporta's[5] second critique is that the exclusion assumption may be

134    violated in this model because "it is possible that temperature affects deforestation not only

135    through malaria, but also through other causal paths," particularly the relationship between

136     temperature and agricultural gross domestic production.[22] In other words, favorable temperatures

137     for mosquitos and malaria parasites may affect deforestation not just through malaria, but by also

138     being favorable agricultural growing conditions, which increase the potential value of forest

139     clearing. We agree that temperature is important to both agriculture and malaria, and that those

140     clearing land may consider the land's growing potential. However, rather than counting the

141     number of days in a 2ºC temperature window during the dry season, we suggest agricultural

142     producers will instead consider the general growing conditions of a region as it relates to

143     commonly grown crops—for example, soil quality, climate, topography, and infrastructure. As

144     land clearing for agriculture is a large and long-term investment, average growing conditions are

145     much more likely to influence clearing decisions than are small deviations in weather from year

146     to year.

147         There are two additional primary reasons that our IV, optimal malaria transmission

148     temperature, is highly unlikely to fail the exclusion criteria. First, we specifically employ

149     municipality "fixed effects" or dummy variables[15] to remove roughly time invariant

150     characteristics specific to each municipality through differencing. Thus, average characteristics

151     (e.g., soil quality, average precipitation, average temperature) that are likely to influence the

152     evolution of regional agricultural land use and the location of processing plants and other

153     infrastructure are removed and the model is identified from deviations from the municipality-

154     specific mean. Second, the range of optimal average temperatures for soybean—Brazil's main

155     crop by area and production[23]—cultivation and development in Brazil is from 20ºC to 35ºC.[24]

156     Recall optimal temperature for malaria transmission is 24ºC to 26ºC, and we use the number of

157     days in the dry season within this narrow temperature band as our instrument. Thus, an

158     additional day at 25ºC relative to 27ºC would be expected to lead to increases in malaria

159　transmission. However, this same change in temperature would likely have a trivial impact on

160　soy yields, as both temperatures are well within the bounds of optimal soy cultivation. Given the

161　breadth of favorable temperatures for soy, it is unlikely that changes in the number of days

162　between 24ºC to 26ºC will influence land clearing decisions for agricultural production.

163　　　We too feel that causal inference approaches hold much promise in disease ecology, and

164　agree that researchers interested in exploring the use of such methods should carefully consider

165　model assumptions. Toward that end, we briefly derive the simplest form of IV to illustrate to

166　potential users what is under the hood of the IV approach and how the exclusion and relevance

167　assumptions function in this technique.

168

169　*Deriving the IV Estimator:* To keep it as intuitive as possible, let us assume a bivariate regression

170　of the form,

171

172
$$y_i = \alpha + \beta x_i + \varepsilon_i \qquad\qquad 1$$

173

174　Where $y_i$ is the outcome variable (e.g., malaria incidence) for observation (e.g., municipality) $i$,

175　$x_i$ is the endogenous explanatory variable (e.g., deforestation), $\varepsilon_i$ is the error term, $\alpha$ is the

176　intercept, and $\beta$ is the coefficient of interest.

177

178　To derive the IV estimator, we can take the covariance of each side of equation 1 with respect to

179　the instrument, $z_i$:

180

181
$$cov(z_i, y_i) = cov(z_i, \alpha) + cov(z_i, \beta x_i) + cov(z_i, \varepsilon_i) \qquad\qquad 2$$

182

$$= 0 + \beta cov(z_i, x_i) + cov(z_i, \varepsilon_i) \qquad\qquad 3$$

184

185  Since $\alpha$ is a constant, and the covariance of a variable with a constant is 0, the first term drops

186  out. Similarly, because $\beta$ is a constant, it can be removed from the covariance. The exclusion

187  assumption of IV is that the instrument $(z_i)$ only affects the outcome through changes in the

188  endogenous variable $(x_i)$, which is more formally written as $cov(z_i, \varepsilon_i) = 0$. Thus with basic

189  rearranging, we have derived the IV estimator $(\beta_{IV})$,

190

$$\beta_{IV} = \frac{cov(z_i, y_i)}{cov(z_i, x_i)} . \qquad\qquad 4$$

192

193  *Consistency of IV:* If we then want to illustrate that the IV estimator is consistent—in other

194  words, as the sample size gets larger and larger the distribution of the estimator converges to the

195  true parameter value—we can plug the right-hand side of equation 1 into $y_i$ in equation 4. We

196  substitute $\beta_{IV}$ with $\widehat{\beta_{IV}}$ since we are considering whether the estimated slope from an IV

197  converges in probability to the true slope $\beta$.

198

$$plim\ \widehat{\beta_{IV}} = \frac{cov(z_i, \alpha + \beta x_i + \varepsilon_i)}{cov(z_i, x_i)} . \qquad\qquad 5$$

200

201  Following a similar logic as with equation 3, equation 5 becomes:

202

203
$$plim \, \widehat{\beta_{IV}} = \frac{\beta cov(z_i, x_i)}{cov(z_i, x_i)} + \frac{cov(z_i, \varepsilon_i)}{cov(z_i, x_i)}.$$
6

204

205    From equation 6, the second assumption of IV becomes evident. The second assumption is the

206    relevance assumption, or that the instrument must be statistically associated with the endogenous

207    variable $(x_i)$. As can be seen in equation 6, this means, in mathematical terms, $cov(z_i, x_i) \neq 0$.

208    Covariance does not imply a direction to the relationship, whether AOD (our instrument)

209    determines deforestation or deforestation determines AOD (or neither) is irrelevant, as it is the

210    covariance between the two that is important.

211

212    By these two assumptions of IV, that $cov(z_i, \varepsilon_i) = 0$ and $cov(z_i, x_i) \neq 0$, equation 6 simplifies

213    to $plim \, \widehat{\beta_{IV}} = \beta$, illustrating IV is a consistent estimator of the true relationship.

214

215    *Weak Instruments:* Equation 6 also illustrates another important aspect when considering the

216    application of instrumental variables, and that is a problem known as "weak instruments." The

217    problem occurs if the exclusion criteria, $cov(z_i, \varepsilon_i) = 0$, fails. Based on the relationship between

218    covariance and correlation (namely, $cov(x, y) = corr(x, y) * \sigma_x \sigma_y$ where $\sigma$ is the standard

219    deviation of each variable) and assuming $cov(z_i, x_i) \neq 0$, we can rewrite equation 6 to illustrate

220    the problem (omitting subscripts for simplicity).

221

222
$$plim \, \widehat{\beta_{IV}} = \beta + \frac{corr(z, \varepsilon) * \sigma_z \sigma_\varepsilon}{corr(z, x) * \sigma_z \sigma_x} = \beta + \frac{corr(z, \varepsilon) * \sigma_\varepsilon}{corr(z, x) * \sigma_x}.$$
7

223

224    If there is a small correlation between the instrument and the error, the last term in equation 7

225    does not drop out and the IV estimator is inconsistent ($plim\ \widehat{\beta_{IV}} \neq \beta$). If $corr(z,\varepsilon)$ is just

226    slightly different from zero and $corr(z,x)$ is much different than zero, the last term is of

227    minimal influence. However, if the instrument is only weakly correlated with the endogenous

228    covariate, the last term of equation 7 can become large. In practice, weak instruments can cause

229    the IV estimator to be severely biased. Since there is no test to validate the exclusion criteria, the

230    strength of the relationship between the instrument and the endogenous variable is very

231    important in practice, and can be formally tested[25] as in the supplementary material from

232    MacDonald and Mordecai.[4]

233

234    *Conclusion:* Understanding the effects of environmental change on infectious disease

235    transmission—from diseases long endemic to the tropics like malaria, to novel emerging

236    pathogens we have yet to discover like SARS-COV-2—is of fundamental and increasing

237    importance. In these complex socio-ecological systems that are difficult to study experimentally,

238    emerging data sources (e.g., high spatio-temporal resolution earth observation data) and causal

239    inference methods (e.g., IV estimation) represent one methodological approach that can help us

240    achieve such clearer understanding.

241

251

252     Author Contact Information:

253     Andrew J. MacDonald: Bren School of Environmental Science and Management, University of

254     California, Santa Barbara, CA 93106-5131; andy.j.macdon@gmail.com

255     Erin A. Mordecai: Department of Biology, Stanford University, Stanford, CA 94305;

256     emordeca@stanford.edu

257

## References

259     1.      Plowright RK, Reaser JK, Locke H, Woodley SJ, Patz JA, Becker DJ, Oppler G, Hudson

260             PJ, Tabor GM, 2021. Land use-induced spillover: a call to action to safeguard

261             environmental, animal, and human health. *Lancet Planet Health* 5(4):e237-e245.

262             doi:10.1016/S2542-5196(21)00031-0.

263     2.      Thomas MB, 2020. Epidemics on the move: Climate change and infectious disease. *PLoS*

264             *Biol* 18(11):e3001013–2. doi:10.1371/journal.pbio.3001013.

265     3.      Jones KE, Patel NG, Levy MA, Storeygard A, Balk D, Gittleman JL, Daszak P, 2008.

266             Global trends in emerging infectious diseases. *Nature* 451(7181):990-993.

267    4.    MacDonald AJ, Mordecai EA, 2019. Amazon deforestation drives malaria transmission,

268          and malaria burden reduces forest clearing. *Proc Natl Acad Sci USA* 116(44):22212-

269          22218. doi:10.1073/pnas.2014828117.

270    5.    Valle D, Laporta GZ, 2021. A Cautionary Tale Regarding the Use of Causal Inference to

271          Study How Environmental Change Influences Tropical Diseases. *Am J Trop Med Hyg*.

272          doi:10.4269/ajtmh.20-1176.

273    6.    Larsen AE, Meng K, Kendall BE, 2019. Causal analysis in control–impact ecological

274          studies with observational data. *Methods Ecol Evol* 10(7):924-934. doi:10.1111/2041-

275          210X.13190.

276    7.    Bonds MH, Dobson AP, Keenan DC, 2012. Disease Ecology, Biodiversity, and the

277          Latitudinal Gradient in Income. *PLoS Biol* 10(12):e1001456.

278          doi:10.1371/journal.pbio.1001456.

279    8.    MacDonald AJ, Larsen AE, Plantinga AJ, 2019. Missing the people for the trees:

280          Identifying coupled natural–human system feedbacks driving the ecology of Lyme

281          disease. *J Appl Ecol* 56(2):354-364. doi:10.1111/1365-2664.13289.

282    9.    Bauhoff S, Busch J, 2020. Does deforestation increase malaria prevalence? Evidence from

283          satellite data and health surveys. *World Dev* 127:104734.

284          doi:10.1016/j.worlddev.2019.104734.

285    10.   Jones IJ, et al., 2020. Improving rural health care reduces illegal logging and conserves

286          carbon in a tropical forest. *Proc Natl Acad Sci USA* 117(45):28515-28524.

287    11.    Garg T, 2019. Ecosystems and human health: The local benefits of forest cover in

288            Indonesia. *J Environ Econ Manage* 98(24):102271. doi:10.1016/j.jeem.2019.102271.

289    12.    Couper LI, MacDonald AJ, Mordecai EA, 2021. Impact of prior and projected climate

290            change on US Lyme disease incidence. *Glob Chang Biol* 27(4):738-754.

291            doi:10.1111/gcb.15435.

292    13.    Larsen AE, MacDonald AJ, Plantinga AJ, 2014. Lyme Disease Risk Influences Human

293            Settlement in the Wildland-Urban Interface: Evidence from a Longitudinal Analysis of

294            Counties in the Northeastern United States. *Am J Trop Med Hyg* 91(4):747-755.

295            doi:10.4269/ajtmh.14-0181.

296    14.    Santos AS, Almeida AN, 2018. The Impact of Deforestation on Malaria Infections in the

297            Brazilian Amazon. *Ecol Econ* 154:247-256.

298    15.    Wooldridge JM. *Econometric Analysis of Cross Section and Panel Data*. first edition.

299            Cambridge, Massachusetts: MIT Press; 2002.

300    16.    Morgan WT, Darbyshire E, Spracklen DV, Artaxo P, Coe H, 2019. Non-deforestation

301            drivers of fires are increasingly important sources of aerosol and carbon dioxide emissions

302            across Amazonia. *Sci Rep* 9:16975. doi:10.1038/s41598-019-53112-6.

303    17.    Aragão LEOC, et al., 2018. 21st Century drought-related fires counteract the decline of

304            Amazon deforestation carbon emissions. *Nat Commun* 9:536. doi:10.1038/s41467-017-

305            02771-y.

306    18.    Chen Y, Morton DC, Jin Y, Collatz G, Kasibhatla PS, van der Werf GR, DeFries RS,

307          Randerson J, 2013. Long-term trends and interannual variability of forest, savanna and

308          agricultural fires in South America. *Carbon Manag* 4(6):617-638. doi:10.4155/cmt.13.61.

309    19.    Hengl T, Wheeler I, 2018. Soil organic carbon content in x 5g / kg at 6 standard depths (0,

310          10, 30, 60, 100 and 200 cm) at 250m resolution (Version v0.2).

311          doi:10.5281/zenodo.2525553.

312    20.    Gorelick N, Hancher M, Dixon M, Ilyushchenko S, Thau D, Moore R, 2017. Google Earth

313          Engine: Planetary-scale geospatial analysis for everyone. *Remote Sens Environ* 202(C):18-

314          27. doi:10.1016/j.rse.2017.06.031.

315    21.    Mordecai EA, et al., 2012. Optimal temperature for malaria transmission is dramatically

316          lower than previously predicted. *Ecol Lett* 16(1):22-30. doi:10.1111/ele.12015.

317    22.    Burke M, Hsiang SM, Miguel E, 2015. Global non-linear effect of temperature on

318          economic production. *Nature* 527(7577):235-239. doi:10.1038/nature15725.

319    23.    Cattelan AJ, Dall'Agnol A, 2018. The rapid soybean growth in Brazil. *OCL* 25(1):D102.

320          doi:10.1051/ocl/2017058.

321    24.    Viana JS, Gonçalves EP, Silva AC, Matos VP, 2013. *Climatic Conditions and Production

322          of Soybean in Northeastern Brazil*. IntechOpen. doi:10.5772/52184.

323    25.    Olea JLM, Pflueger C, 2013. A robust test for weak instruments. *J Bus Econ Stat*

324          31(3):358-369. doi:10.1080/00401706.2013.806694.

325