

Review of « Weather and air pollution influences on solar energy performance in West Africa: A Bayesian nonlinear mixed-effects approach », by Konin Pierre-Claver Kakou et al. Review by Quentin Libois.

### **General comments**

This paper investigates the dependence of Global Horizontal Irradiance (GHI) on weather conditions at the surface (temperature, humidity, wind speed), cloud cover and aerosol concentration, for 12 stations located in Côte d'Ivoire. To this end it relies on a sophisticated Bayesian model to derive the posterior distributions of the parameters relating the predictors to the predicted variable. The paper is overall well written, in particular the introduction that gives a nice overview of the related works, with a special focusing on the uncertainty estimation for the prediction models.

The paper is very technical, particularly hard to follow for non-statisticians, and lacks of physical interpretations. In many places reordering would help, and some paragraphs are not specifically useful and could be removed or at least condensed. The section dedicated to the analysis of the results is very general and does not demonstrate the interest of the work in terms of processes understanding or improved prediction capability, which somehow deserves the primary objective of the paper. Likewise the discussion is very limited, while obviously there would be much to say. Unless the paper is augmented in terms of original content and its added value to the field is better emphasized, I doubt it can be of great interest for the readers of GMD, except maybe for the pure methodological approach that could be used by others.

### **Specific comments**

1) There has been a lot of recent work on the forecasting of solar energy. It should be better emphasized what is the originality of the present study. In terms of methodology, efficiency (then it should be compared to state-of-the-art alternative models), location, etc. These elements should be motivated in the introduction, and highlighted in the conclusion. Likewise, in Section 3.3 some dependences between GHI and meteorological variables are highlighted, that could probably be mentioned in a dedicated paragraph of the introduction, to point out what is already known in terms of such correlations.

2) It is not clear why it is useful to derive such prediction models, instead of simply measuring the GHI. Maybe the point would be to use forecasts of meteorological parameters and to predict GHI from that. In this case, would the estimation be better than the GHI forecasts itself (from satellite products of weather models?). A comparison could help convince the reader of the utility of the new model. Otherwise, the physical understanding of the links can in itself be valuable, but it is very weak in this paper, so it does not bring much novel physical insight to the community. The applications could be mentioned in the discussion, or at least in the opening of the conclusion.

3) The paper is very focused on the statistical model, in a way that is hard to follow for non-experts. To increase the readability it would be great to better explain the physics behind the parameters, scores, tools etc. This probably implies a partial rewriting of the Method section.

4) As said at point 2) the physical interpretation of the results is extremely limited. While the method only puts forward correlations, the authors try to find some causality, and for each explanation we could argue for the opposite causal chain. Also the fact that covariance between variables is not really discussed is such that apparent correlations between a variable and GHI can be directly related to the correlation of one variable with the other. The opposite may also happen when the authors do not see a strong impact of PM<sub>2.5</sub>, the inclusion of which sounded however as

an originality of the paper. In practice, the authors mention many physical processes that can explain the link between meteorological variables and GHI, but they struggle to convince the reader about what really drives the correlations. Section 3.3 hence appears as a sum of physical explanations and suggestions without clear, ordered conclusions. This is critical for a method that puts forward the interpretability.

5) The paper is somehow too short, with the largest share dedicated to the statistical model. It would benefit from additional figures (ex: an illustration of the performances of the final model in terms of predictions vs observations). The discussion is definitely too short as well, and very weak as is.

### **Technical comments**

l.5: “this relationship” is poorly defined as so far the predictors are not explicated. Do they correspond to individual weather variables such as temperature, humidity, cloud cover, liquid water path? In which case it would be obvious that solar radiation is not linearly related to some of those variables.

l.8: “prior knowledge” of what?

l.12: if the graph neural networks is a perspective (is it an alternative to the Bayesian approach?), it should not appear (at least that early) in the abstract.

l.16: as the domain of application (space and time) has not be defined, allusion to “similar tropical climates” is unclear. The abstract clearly lacks of a mention of the data used.

l.26: it is not clear weather “predict” refers to some leadtime (in the future) or to estimate the present solar radiation from present weather characteristics

l.46: when citing a paper its content should be a bit more detailed. Here the reported result is too general (no information on location, type of situations etc.).

l.46-47: remove parentheses for citations. Beware also with the use of `\cite{}`

l.50: it’s not clear why adding a possibly useless predictor can be detrimental to the model performance

l.51: when referring to air pollutants, it would be more rigorous to mention the actual quantities used (concentrations near the surface, total column, AOD?). Also we expect given this sentence that you’ll use BC in your study.

l.68: “efficiency” has not been defined, and can be misinterpreted

l.70: in the mentioned “products” the variables are not really “parameterized”, they’re rather “estimated”

l.70-71: “due to limited data availability and an incomplete understanding of their regional dynamics” is unclear. You can mention their intrinsic limitations (spatial resolution, precision, etc.)

l.72: satellite products of which quantities?

1.75: why would satellite miss these low-level clouds? Small size, poor contrast with the surface? Any reference to support this statement?

1.75: is CAMS a satellite product? I mostly know the CAMS model

1.75: “the authors of the present study found” could be turned into “it was recently shown”, as the citation comes in the end

1.76: please specify the considered variables

1.84: an outline of the paper at the end of the introduction would be welcome. To detail what kind of data are used in particular.

1.87: could the twelve stations be highlighted in Fig. 1?

1.109: the objective **of** this study

1.113: in details

1.113: the references should not be in parentheses

1.114: what is cloud cover here? A fraction of cloud (%) or a cloud mask at the spatial resolution of the satellite product ? In this paragraph the satellite products used and their characteristics (resolution, frequency, type (column, profile), etc.) should be detailed. Also, do not mention satellite products if a reanalysis is actually used (for aerosols).

1.119: any motivation to only consider PM<sub>2.5</sub> (not larger particles) to estimate the impact on solar radiation? Also do you mean here mass concentration again?

Eq.1: I have the feeling that if some particles grow they become larger than 2.5 micrometers, hence would not be included anymore in [PM<sub>2.5</sub>]. Does it deserve some explanation?

1.125: what is a wet/dry relative humidity? 0 and 100%? Is this naming standard? As it sounds a bit awkward

1.128: should this ratio f<sub>OM</sub>:OC appear in Eq. 1?

1.146: “pollution” seems loosely defined. Does it correspond to all aerosols? Why not using aerosols instead?

1.153: not clear what you mean by “attenuation” here. The fact that the air mass increases as the Sun is low, or just the fact that the Sun is low hence the GHI lower?

1.155: the wording is surprising. Maybe just state that all values outside the range were discarded because unexpected? But then what does it mean if you discard some values (and how many?), that the measurements are not reliable? Hence they could be unreliable also for values within the range?

1.157: could you specify for which variables the outliers are removed?

1.164: I think you can use GHI here

1.166: it seems that some predictors are missing

Table 2: the acronyms should be understandable

l.182: not clear in Eq. 4 what is known and what should be retrieved (the parameters)

l.192: what are the “global parameters”?

l.192: “combination” is unclear. Is it a sum, a product? How are the actual nonlinear functions chosen for each predictor?

l.194: even after normalization there remains an impact of the diurnal cycle due to the air mass

l.203: for non-experts, maybe explain why values larger than 1 can be observed, otherwise it suggests that the measurements could be erroneous

l.205: should this Beta distribution be explicit? At least I’m not familiar with it

l.207: not clear where the precision parameter comes from, and what value it takes

l.225: maybe explicit that the sigmoid fit was always the best one, except for  $P_{it}$ . Also not clear how you compare R values for different sites. Do you choose the function that has the maximum mean R value?

l.228: should it be  $H = 3$ ? Also define  $d_i$

Figure2: the caption contains important information about the ranges of observations used, that may deserve to be in the main text. Also, what does it mean to obtain fits for different ranges, sometimes inconsistent? In terms of general applicability of the final regression? This point probably deserves more explanations

l.231: again it is not clear what the BNLME model does or finds once you’ve found the regression functions for each variable. What are the free parameters to be optimized?

l.233: not clear what the “raw” covariates are

l.235: what is the motivation of using so many models, while one would expect only the more detailed one to perform best (according to Eq. 8)

l.254: the references do not appear correctly here (should be in parentheses)

l.260: it is not clear so far why you need to define such sky conditions. Maybe add a sentence to tell why you use such a classification

Section 3.1. Not clear what is the added value of this very descriptive paragraph, given that Fig. 2 contains already much information about the variables. In the text details are provided for some arbitrary stations, not all, with no specific motivation for the choice made. There does not seem to be critical results highlighted in that paragraph, I believe Table 4 could be self-sufficient. Saying that the minimum wind speed is  $0 \text{ m s}^{-1}$  and the minimum insolation nearly  $0 \text{ W m}^{-2}$  as well is not very useful.

l.269: Highest variability **for** this variable

l.291: I confess that `elpd.loo` is not very meaningful to me. I wonder if it would require more low-level explanation. Likewise, what is the model “weight”?

Figure 4: what are ELPD and ELPD difference? And how is it meant to be interpreted?

l.301: MCMC has not been defined, neither the corresponding diagnostics

l.304: I’m not sure I understand what parameter is the intercept

Fig. 5: to what station correspond each color in the first row? Also, there are different lines. To what do they correspond?

l.309: does  $b_i$  correspond to the first row (beta)? In that case this is misleading, as beta appears for each variable as well in Eq. 8. Also the variable  $b$  appears twice in Eq. 8 (at the beginning and for the temporal variations).

l.324: I’m not sure the prior distributions of all parameters have been introduced

Table 7: It is not clear to me whether there should be a set of “parameters distributions” for the whole country, or one for each station. In the latter what is the mean in Table 7? The mean of all the posterior distributions? It’s hard to understand how a set of posterior distributions is treated and why it would be meaningful to give these statistics (compared to showing all the posterior distributions).

l.349-350: this point could have been clarified earlier

l.354-357: I’m not convinced by this link between surface temperature and cloudiness, and by the dissipation argument. Could it be reversed? Less clouds means more GHI hence higher temperatures? Likewise the impact of wind is very questionable. Wind could suspend dust and decrease GHI, and often wind is associated with cloudy conditions. You’re looking at correlations hence it is very difficult to find causal effects, which could be one way or the other.

l.365: what do you mean by “attenuation effect of relative humidity”? The actual absorption by water vapor (to be analyzed in clear-sky conditions only) or the fact that the dependence is stronger wrt humidity (through the intermediate formation of clouds maybe)?

l.375: again, there might be strong covariance that could mask the first-order effect of pollution, for instance if pollution is correlated with wind

l.384: it is the only place where sky conditions are used, and this is marginal. Would it be useful to show a Figure to illustrate this particular case?

l.391: indeed AOD would be more relevant, and I see no particular reason not to use such a product, for instance from AERUS-GEO (<https://www.icare.univ-lille.fr/projects/user-driven-projects/aerus-geo/>)