

General comments

C: *This paper explores whether satellite and reanalysis data for rainfall and soil moisture can be combined using machine learning methods to assess, in an objective way, whether floods or droughts are happening or have recently occurred. This is placed in the context of improving index insurance. The paper is extremely detailed in terms of how the machine learning models are constructed, and validation metrics.*

R: Dear Reviewer,

Thank you very much for your time and effort reviewing our manuscript. This response (R) carefully addresses all the comments (C). Where deemed appropriate, modifications to the manuscript are proposed (red underlined text indicates additions to the manuscript, blue strikethrough text indicates removed text).

C: *My main comment is: the paper is very heavy on text-book style review of methods (which isn't a bad thing), and very heavy on technical detail (which isn't a bad thing), but lacks any exhibits that show clearly whether the methods actually work or not. There are masses of technical validation metrics. But what I personally would like to see are some results along the lines of: a) we took the data shown in figure 7 (predicting this data is what the whole thing is about in the end) b) we split that data in half, trained the models on one half, chose the best model, and tested it on the other half c) and for the single best model, here's a picture that shows the results of that side by side with the actual floods and droughts that occurred in the validation period. Did it capture them all, or half of them, or none of them? d) then I'd be able to look at that and make a judgement as to whether the method works or not.*

R:

After analysing this comment together with the rest of the review, our understanding is that the Reviewer has some concerns regarding the performances of the models, which originate in part from how he/she perceives that the validation process was carried out. Upon critically reviewing the original manuscript, we believe that this is likely because the validation process was not sufficiently well explained. More specifically, Section 2.3.3 describes in detail which are the best practices used in the validation of a machine learning model without stating clearly what was done in our work. Hence, we suggest the following change at line 318 when we describe what was done in our work regarding the splitting and validation, which we hope clarifies this aspect.

“In this work, the proper training of the NN was exerted splitting the dataset in 3 parts: training (60%), validation (15%) and testing (25%) set. During training, the neural network used only the training set, evaluating the loss on the validation set at each iteration of the training process. After the training, the performance of the model was evaluated on the testing set that the model has never seen. Concerning the SVM, a k-fold cross validation (Mosteller and Tukey, 1968) was used to validate the SVM model, using 5 folds created by preserving the percentage of

sample of each class, the algorithm was therefore trained on 80% of the data and its performances were evaluated on 20% of the remaining data that the model has never seen whilst for NN, TensorFlow allows the user to declare a percentage of the data that is retained as validation data at each iteration of the training loop, therefore, embedding the validation process into the construction of the model.”

We hope that the above changes clarify how the splitting of the dataset and the training process were performed, which is a central aspect for the construction of the models and the robust evaluation of their performances.

Going back to the Reviewer’s comment more specifically, we would like to highlight that the objective of verification metrics such as the ones we used in our manuscript is summarizing the overall prediction quality of a set of predictions - which is the same as to “assess whether the methods actually work or not”. Notwithstanding, we do recognize that looking at an actual plot of observations versus models predictions can sometimes be an easier way to have a feeling for how the model is performing and therefore we considered adding such a plot to the article. Below is a figure depicting, for the days and weeks belonging only to the testing set (i.e., data that the models have never seen before), observed flood and drought events and the corresponding predictions from the three methods discussed in the manuscript.

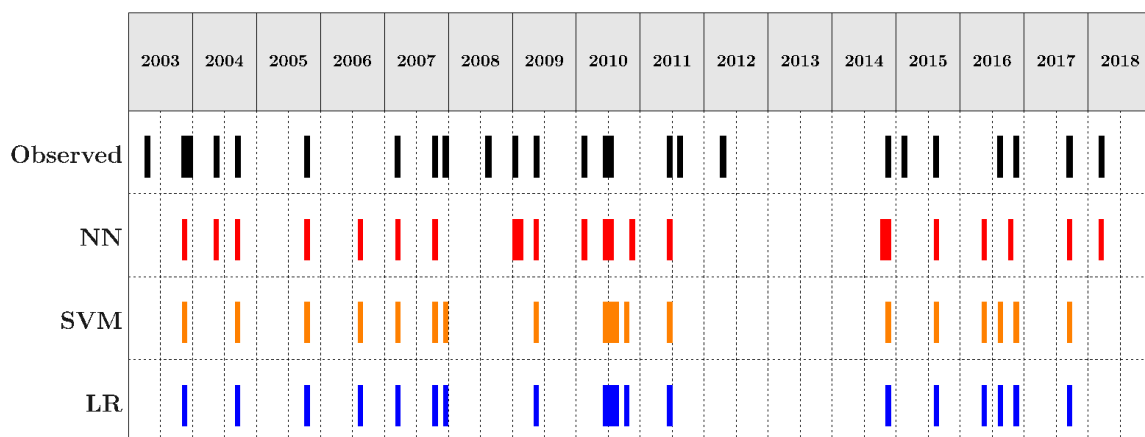


Figure 1: Comparison of prediction over the testing set of the three methods. Flood Case.

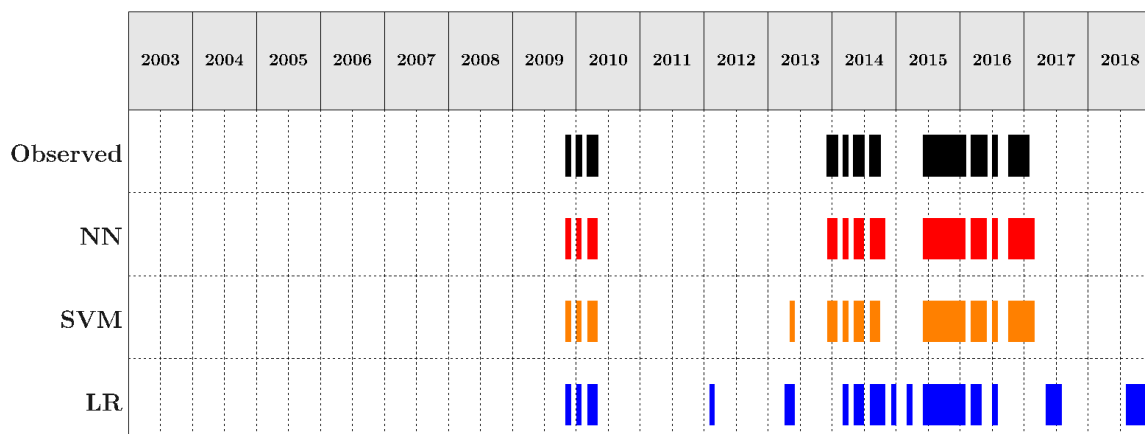


Figure 2: Comparison of prediction over the testing set of the three methods. Drought Case.

The above plot is essentially a visual representation of what the metrics reported in Table 7 and 8 report. The question is whether adding such plots would objectively improve the article. After careful consideration, we believe that while plots such as these may be more visually appealing, they do not significantly contribute toward a better interpretation of the results. Therefore, we propose to keep the numerical representation, which we believe provides an easier way to compare results among methods and is able to summarize in a short and quantitative way the results, while including the above plots as Supplementary material.

Specific comments

C: There's a whole discussion about training and validation data, but then in the end it's not clear how the data is actually split into training and validation data (relates to point 1 above), in relation to Figure 7. The construction of the validation is critical for us to be able to understand whether there's anything in this or not, especially since a large part of the scientific community associates the word 'machine learning' with 'overfitting', and will be sceptical.

R: We agree that proper validation is crucial and paid careful attention to its construction. We hope the changes suggested above provide a good improvement to the reading experience and dispel any doubt about the partitioning of the dataset.

C: With such a small amount of data, and after testing so many models and configurations (line 341: 'almost boundless domain of model configurations'), it seems to me that overfitting is quite likely. a) Could the authors elaborate on why testing so many configurations doesn't lead to overfitting? b) And if you are evaluating the models against each other using the validation dataset, of course one model will do best. How do we know that the model that does best would genuinely do best in a true out of sample sense? Don't you need another level of cross-validation?

R: I'll address the reviewer questions in order:

- a) The different model configurations tested do not share the training process. Each model, initialized with different parameters and configurations the way we presented in the manuscript, has its own training process, hence, the training of one configuration does not influence the others. We would like to emphasize that having different model configurations to train is not to overfitting, which by definition refers to the ability of a model to reproduce predictions or analysis valid only for a particular set of data. This is prevented in our work by providing an accurate splitting of the data, as described above, and by taking other measures during the training of the models that are recalled later on in this response.
- b) The model performances reported were obtained from the testing partition of the data that we consider a "true out of sample". The definition of the testing dataset will be

improved and clarified with the text addition at line 318 mentioned above. This concept is recalled in the result section at line 473.

C: *Line 18 says \$3.3B. This is wrong by several orders of magnitude. Individual events during that period were in excess of \$50B (since at this point you are talking globally).*

R: We will change the number to the correct amount, which is \$3300B. (Hoeppel, 2016)

C: *The word 'loss' is used with two different meanings, as far as I can tell. Line 105=loss in the usual sense of damages, vs line 249 in a technical sense. This is a bit confusing. Different terminology should be used, somehow, to avoid this.*

R: We understand that the usage of loss might create some confusion, but we also acknowledge that the term loss in these two contexts has a specific meaning that would be lost changing the terminology. Therefore, we suggest introducing the acronym LF defining the loss function throughout the rest of the manuscript.

C: *I think it should be made clear that the runoff model – flood intensity relationships are simplistic relative to state of the art runoff and flood modelling as practised by hydrologists*

R: Complex state-of-the-art models are typically not considered in parametric insurance products, as such products are meant to be based on simple indices obtained directly from environmental variables without the need for large modelling efforts that typically require much larger amounts of data. Even if this is somewhat implied, we agree that the simplified approach should be pointed out and therefore propose the following addition to highlight this aspect without going into the details of the runoff/flood modelling practices, which we believe are beyond the scope of the manuscript:

At line 157:

“To achieve this, we adopt a variable transformation to emulate, in a simplified manner, the physical processes behind the occurrence of flood damage due to rainfall, (...).”

C: *Line 176 refers to loss data. What is this loss data?*

R: In this instance we refer to reported occurrence of loss data. We understand the writing might be a bit troublesome, thus, we suggest the following change at line 176 to avoid any misunderstanding, recalling terminology introduced before in the manuscript and used throughout the whole paper:

“...fitting a logistic regression model to concurrent potential flood intensity and reported occurrences of losses caused by flood events data, and maximizing the likelihood using a quasi-Newton method.”

C: *Line 319, there is a comment that TensorFlow allows ‘embedding the validation process into the construction of the model’. That sounds like overfitting to me. Please explain how this is consistent with the claim that the data is really being split in order to do out of sample validation.*

R: We agree that this statement may be somewhat unclear and have proposed to replace it, as described above. Furthermore, we have made several proposals for improvement regarding model validation, which we hope clarify this topic.

Nevertheless, we would still like to provide information about the statement originally reported in the manuscript. When training the model, the function that is called in the code to carry out the training allows, among its parameters, to indicate the way you want to perform the validation. You can either directly pass the validation dataset to the function or declare the percentage of data you want to retain from the training set as part of the validation. In both cases, the validation datasets are not used during the training of the model but to evaluate the loss function during the training.

C: *Is reanalysis data really available soon enough to be useful? I thought it usually appears at least a year or two later, but maybe I’m wrong.*

R: The reanalysis data we used in our work are updated daily with a latency of about 5 days, as reported in table 2 and by the Copernicus documentation ([Copernicus ERA-5](#))

C: *There should be a bit more discussion about the problems with satellite data and reanalyses (i.e., talk about the reasons why these datasets aren’t really used at present for index insurance purposes, even after 20 years of academics suggesting that they should be).*

R: While some limitations are discussed in the introduction, we recognize that highlighting more about the shortcomings of using this type of data could be beneficial to the manuscript.

Addition to the manuscript at line 48:

“Satellite images are often available with high spatial resolution, but records are still short, with a maximum duration of around 30 years. Reanalysis, on the other hand, provides longer time series but tends to have a coarser spatial resolution. Moreover, satellite data should be checked for consistency with ground measurement which is not always feasible when the network of ground instruments is inadequate or non-existent (Loew et al., 2017). Although using satellite data has its own limitations, various index-based insurance products, exploiting remote-sensing data and reanalysis, have been developed in data sparse regions such as Africa and Latin America (Awondo, 2018; African Union, 2021;The World Bank, 2008). The combined use of various sources of information to detect the occurrence of extreme events is

valuable, since it can significantly improve the ability to correctly detect extreme events (Chiang et al., 2007) and a proper index design helps addressing the limitations brought by satellite data, as underlined in Black (2016).

C: As far as I understand it, there has been no comparison here with standard methods for assessing whether an event has occurred, which are based on rain gauges, levels of river flow, etc. That should be pointed out.

R: The historical event catalogue contains information about whether events have occurred or not. Standard methods, which we presume refer to threshold-based methods based on rain or stream gauges, are not applicable in this case as the availability of such data in the Dominican context is scarce. Nevertheless, note that in parametric insurance the use of such data is not optimal, as described in the Introduction.

C: Are there any further diagnostics that could be produced to help show that the model is really doing something sensible, to help allay the suspicion that some readers may have that it's all just over-fitted

R: We hope the discussion and clarification throughout this response have alleviated some of the concerns the review raised regarding overfitting. Nonetheless, we would like to summarize the measures taken in our work to prevent overfitting that are treated in the paper:

- At line 253, we discuss the role of monitoring the training and validation in avoid overfitting and how countermeasures to stop the training can be taken whereas is needed. In practice, the training is stopped with a TensorFlow's call back called "EarlyStopping" that stops the training when a monitored metric, in our case the loss function, has stopped improving. We mentioned in the manuscript how this lack of improvements might very well be unhealthy for the model therefore granting the stop of the training.
- The second paramount action taken against overfitting is the splitting of the data. We clarified how the splitting was done in the previous answers, here we would like to reiterate that the results presented were obtained from data the the model had never seen.

Technical corrections

C: you say T_t , but don't you mean Y_t ?

R: Checked and corrected

C: SPI, 3 etc need to be defined. I can guess what they are, but they should be defined.

R: We agree with the reviewer that the terminology used should be explained clearly. We suggest the following addition to the manuscript:

“In this study SPI1, SPI3, SPI6 and SPI12 were computed, where the numeric values in the acronym refer to the period of accumulation in months (e.g. SPI3 indicates the standard precipitation index computed over a three months accumulation period).”

C: *is that citation really correct? Is the person's name just M?*

R: Will be corrected in “(Hossin and Sulainman, 2015)”

C: *i.e. and e.g. are usually followed by commas I believe*

R: Will be corrected where needed

C: *the plural of reanalysis is reanalyses*

R: Will be corrected where needed

References

- African Union. (2021). *African Risk Capacity: Transforming disaster risk management & financing in Africa*. <https://www.africanriskcapacity.org/>
- Awondo, S. N. (2018). SC. *Journal of Development Economics*. <https://doi.org/10.1016/j.jdeveco.2018.10.004>
- Loew, A., Bell, W., Brocca, L., Bulgin, C. E., Burdanowitz, J., Calbet, X., Donner, R. V., Ghent, D., Gruber, A., Kaminski, T., Kinzel, J., Klepp, C., Lambert, J. C., Schaepman-Strub, G., Schröder, M., & Verhoelst, T. (2017). Validation practices for satellite-based Earth observation data across communities. *Reviews of Geophysics*, 55(3), 779–817. <https://doi.org/10.1002/2017RG000562>
- The World Bank. (2008). *Operational Innovations: Providing Immediate Funding After Natural Disasters*.