



Transitioning Machine Learning from Theory to Practice in Natural Resources Management

Sheila M. Saia^a, Natalie Nelson^{a,*}, Anders S. Huseth^b, Khara Grieger^c, Brian J. Reich^d

^a Biological and Agricultural Engineering, North Carolina State University, Raleigh, NC

^b Entomology and Plant Pathology, North Carolina State University, Raleigh, NC

^c Applied Ecology, North Carolina State University, Raleigh, NC

^d Statistics, North Carolina State University, Raleigh, NC

ARTICLE INFO

Keywords:

Machine learning
Natural resources management
Stakeholders
Decision-support tools
Decision-making
Process-based modeling

what is process-based
modeling?

ABSTRACT

Advances in sensing and computation have accelerated at unprecedented rates and scales, in turn creating new opportunities for natural resources managers to improve adaptive and predictive management practices by coupling large environmental datasets with machine learning (ML). Yet, to date, ML models often remain inaccessible to managers working outside of academic research. To identify challenges preventing natural resources managers from putting ML into practice more broadly, we convened a group of 23 stakeholders (i.e., applied researchers and practitioners) who model and analyze data collected from environmental and agricultural systems. Workshop participants shared many barriers regarding their perceptions of, and experiences with, ML modeling. These barriers emphasized three main areas of concern: ML model transparency, availability of educational resources, and the role of process-based understanding in ML model development. Informed by workshop participant input, we offer recommendations on how the ecological modelling community can overcome key barriers preventing ML model use in natural resources management and advance the profession towards data-driven decision-making.

1. From promise to practice

“Machine learning” (ML) describes a class of algorithms that do not need to be explicitly programmed *a priori* and are highly effective at learning, and making predictions from, patterns in data (Goodfellow et al., 2016; LeCun et al., 2015; Thessen, 2016). Because these approaches are skilled at predicting complex responses from diverse data types, ML is increasingly relevant in the modern era, especially when advances in sensing and computation allow for the natural world to be observed at extraordinary rates and scales (Farley et al., 2018; Lausch et al., 2015; Rode et al., 2016). Despite overlap between ML models and classical statistical models, the motivations for applying these approaches differ. ML models typically focus on prediction, whereas classical statistical models emphasize hypothesis testing and uncertainty quantification (Breiman, 2001; Donoho, 2017). As a result of these differences in motivation, ML models are well-suited to predict nuanced and nonlinear relationships from large, high-resolution datasets (Olden et al., 2008) while classical statistical models (e.g., linear regression) are well-suited to maximize information from small, carefully curated datasets (Hampton et al., 2013). As our capacity to

observe the environment and use these observations for prediction grows, so will the role of ML models in natural resources management.

Leading scientific organizations have promoted the promise of ML models to advance natural resources management by uncovering patterns in large and diverse environmental datasets, and leveraging these relationships to expand and enhance predictive modeling capacity (NASEM, 2019, 2018; WEF, 2018). For example, the World Economic Forum's 2018 report on *Harnessing Artificial Intelligence for the Earth* describes artificial intelligence as key for developing solutions to wide ranging societal challenges such as water availability, food security, and biodiversity conservation (WEF, 2018). Yet, despite growing excitement about artificial intelligence and data science, applying ML models to explore environmental data and develop predictive decision-support tools remains a significant challenge for practitioners working in the natural sciences. Reported barriers to the use of ML models include data-specific challenges (e.g., bias, heterogeneity, size, missing observations), poor accessibility to computational tools and training, and limited knowledge transfer between data scientists, environmental scientists, natural resources managers, and policymakers (Faghmous and Kumar, 2014; Hampton et al., 2017; Kamilaris et al., 2017;

* Correspondence to: Mailing address: Campus Box 7625, Raleigh, NC, 27695.

E-mail address: nnelson4@ncsu.edu (N. Nelson).

<https://doi.org/10.1016/j.ecolmodel.2020.109257>

Thessen, 2016). Although the literature summarizes technical and training challenges hindering the adoption of ML models outside of the computational sciences (e.g., lack of interdisciplinary collaboration; Wagstaff, 2012), few articles offer specific recommendations for actions that may facilitate meaningful and responsible implementation of ML models for decision-making in natural resources contexts.

In an effort to contribute meaningful guidance as to how researchers may increase the adoption of ML models in practice, we invited a group of 23 natural resource management practitioners and researchers to engage in a one-day, face-to-face stakeholder workshop in February 2020, held at North Carolina State University in Raleigh, North Carolina (NC), USA. We invited stakeholders who represented a wide range of intersecting values, knowledge of ML models, sector expertise (i.e., water management, crop production, aquaculture, animal agriculture, air quality, and forestry), and organizations (i.e., federal and local government agencies, multinational companies, engineering consultancies, academia, cooperative extension). The stakeholder workshop was intended for preliminary information gathering (see workshop discussion questions in Table S1). The workshop was not intended to represent a statistically-significant group of stakeholders interested in using ML models for natural resources management. After the workshop, we synthesized responses and feedback from workshop participants and identified three key categories of barriers to ML model adoption: communication, educational resources, and synergies with process-based models. Based on these findings, we provide three recommendations for researchers who are considering using ML models or facilitating the use of ML models for natural resources management in practice. While the stakeholder workshop does not represent a statistically-significant group of stakeholders, we believe our key findings are nonetheless beneficial to researchers involved in applying ML models to natural resources management and communicating ML model results to decision makers.

2. Recommendations to improve ml adoption

1. Improve ML transparency and avoid framing ML models as “black boxes”

Workshop participants expressed concerns that ML models may often be perceived as opaque and inscrutable, thereby preventing their use in practical decision making (e.g., public safety planning, regulatory agency permitting). More specifically, researchers often refer to ML models as “black boxes” because their structures and learned relationships are not as readily interpretable as differential equations and classical statistical models. Workshop participants also viewed the difficulties of interpreting ML model results as being further complicated by the current lack of consensus surrounding the definition and scope of ML. The overlap between ML modeling and classical statistical modeling was confusing to those outside the computational sciences. Without clear, consistent, and easy-to-understand descriptions of ML model structure and scope, stakeholders may view these approaches as too uncertain or risky for use as decision-support tools in natural resource management.

Given workshop participants' concerns about the potential for ML modeling to have ill-defined scope and produce results that are difficult to interpret, we recommend the development of guidelines that work towards improving consensus in scientific messaging on the definition and scope of ML while also revisiting narratives that position ML models as “black boxes”. Descriptions of ML models as “black boxes” implies limited understanding of how their underlying algorithms operate. Though inspecting the inner workings of ML models requires additional effort, researchers, including those outside of computer and statistical sciences, have developed useful and effective approaches for examining ML models and casting light on their internal structures. For example, the *Exploratory Data Analysis using Random Forests* (edarf) R package (<https://github.com/zmjones/edarf>), developed by political scientists, includes

functions to explore features of random forest models such as predictor variable importance and partial relationships between predictor and response variables (Jones and Linder, 2015, 2016). Similarly, the Connection Weights Approach to estimating predictor variable importance (Olden et al., 2004) and *NeuralNetTools* R package (Beck, 2018), developed by a conservation biologist, both facilitate interpretation of supervised neural network models. Additionally, posterior analysis of ML model predictions using interpretation algorithms such as Shapley values (Lundberg et al., 2020) or local interpretable model-agnostic explanations (LIME; Ribeiro et al., 2016) may improve trust in model outputs. However, not all ML model architectures are easy to explore. For example, deep neural networks, which have hundreds or thousands of middle layers, also referred to as “hidden layers” (LeCun et al., 2015; Shen, 2018), are more difficult to interpret compared to simpler ML models with only a one or two middle layers (e.g., multilayer perceptron neural networks). Continued advancement in tools that expose the inner workings of ML models may help improve trust in model predictions, thereby increasing the value of ML models for natural resources management research and practice.

Open and participatory science practices that foster information transparency and co-development of research priorities between researchers and stakeholders may also help address concerns regarding ML models transparency. When applied across the entire research process (i.e., from formation of research question to publication of data and research findings), these practices strive to generate research products that are more inclusive, effective, transparent, reproducible, and discoverable to researchers and stakeholders (Bartling and Friesike, 2014; Hampton et al., 2015; Lowndes et al., 2017; Norström et al., 2020; Woelfle et al., 2011).

2. Develop educational resources on the use of ML models, including descriptive case studies from real-world contexts

Workshop participants emphasized the need for educational materials and case studies on ML modeling that were relevant to natural resources management. While most workshop participants were aware of ML models, many were overwhelmed by the range of ML modeling options, dataset sizes, and computing needs. They asked for specific guidelines and training on technical topics including: data discovery and cleaning, data quality assurance and control, appropriate data requirements (e.g., temporal duration, percent dataset completeness), trusted open-source ML modeling tools, criteria for selecting between various ML modeling approaches and advanced computing resources (e.g., in the form of flow charts), setting-up ML models to be run “in production”, interpretation of ML model outputs and model uncertainty, and limitations of ML modeling. They also asked for guidance on non-technical subjects, including what ethical considerations (e.g., data ownership and privacy, checking for model biases) to make when using ML models for prediction purposes, as well as how to communicate results to various levels of decision-makers, from the general public to elected officials and company leadership.

Workshop participants had many recommendations for how researchers could improve educational resources and accessibility of ML modeling approaches. In particular, workshop participants advocated for the development of case studies that were easy to follow and included model training, tuning, and testing protocols for non-experts making decisions at various spatial scales (e.g., field, region) and time scales (e.g., short-term/emergency, long-term planning). Their suggestion to develop case studies was made in light of the fact that many scientific articles presenting ML modeling applications in the natural sciences are written for ML model experts rather than new users. Therefore, we recommend researchers publishing ML modeling studies relevant to natural resources management consider expanding methods sections and/or supplementary materials to include summaries that contextualize, justify, and describe the use of ML modeling approaches in a way that is well suited for

new ML modelers. Additionally, case studies that provide guidance on how best to translate ML model architectures and outputs for decision-makers may be particularly helpful in improving ML adoption among practitioners.

Currently, many examples demonstrating ML model training, tuning, and validation are presented in the context of software tools (e.g., R package vignettes); however, there is an opportunity to develop ML-specific case studies that go beyond software tool development to improve communication and education strategies. Specifically, these strategies may help bridge the gaps between model predictions, model interpretations, and informed management decisions. Importantly, the co-development of case studies and other educational materials by stakeholders and researchers is needed to ensure these materials meet the needs and interests of stakeholders.

3. Provide guidance on how and when process-based understanding informs ML model architecture

Given the widespread use and trust in established natural resources management methods that rely on process-based models, workshop participants expected to encounter resistance from support staff, leadership, and decision-makers when initiating conversations about adopting ML models for natural resources management. They explained that this resistance likely stems from several barriers. First, workshop participants perceived new methods like ML models as more uncertain than process-based modeling standards, which are regarded as trusted decision-support tools because they encapsulate current knowledge and expertise on underlying processes driving ecological systems (Faticchi et al., 2016; Hipsey et al., 2015; NRC, 2007; Robson et al., 2008). Second, workshop participants noted their unfamiliarity with implementing ML modeling (see Recommendation #2). Last, they were concerned that ML model results may be difficult to interpret (see Recommendation #1) or hinge on spurious relationships in the data that do not uphold process-based understanding of ecological system dynamics.

Considering frequent preferences for process-based models and workshop participants' concerns with ML models, we recommend the development of clear and easy-to-follow guidelines on how non-expert ML modelers can use their knowledge of process-based models to inform ML model development for natural resources management. Applications that bridge ML modeling and process-based modeling, such as theory- or process-guided ML modeling (Faghmous and Kumar, 2014; Hanson et al., 2020; Karpatne et al., 2017; Read et al., 2019), present ML modeling in intuitive and defensible ways for model practitioners. Moreover, ML models are well suited to address limitations of process-based models, such as reducing uncertainty in process-based model parameter estimates (e.g., Gentine et al., 2018) and improving process-based model prediction accuracy (e.g., Read et al., 2019). ML models may help identify novel patterns in environmental data, establish new working hypotheses of underlying mechanisms, and facilitate new field and process-based model experiments to test these hypotheses (Peters et al., 2014; Shen, 2018; Shen et al., 2018). Thus, when developing guidance and case studies demonstrating the utility and value of ML models (i.e., Recommendation #2), we recommend researchers describe how process-level understanding influenced their ML modeling workflows and present ML models as complementary, not contradictory, to process-based models. Last, researchers may consider engaging in participatory research to address how process understanding informs ML model workflows (Norström et al., 2020). In this case, participatory research may lead to co-production of new modeling approaches and model-derived insights.

3. Closing Remarks

As researchers and professionals in the natural sciences apply innovative ML models to manage natural resources in increasingly diverse

disciplines, a firm understanding of the goals, ethics, and interpretations of analytical outcomes will be essential. While our stakeholder workshop was designed for preliminary information gathering, we synthesized and shared important findings from the workshop to provide guidance and recommendations on how improvements in the field of ML can accelerate adoption of ML models for natural resources management. We call on researchers who already work at the intersection of environmental and data sciences to support initiatives that translate the utility of ML approaches to practitioners and, ultimately, advance predictive and adaptive management of natural resources through ML model applications.

Acknowledgements

The authors do not have any conflicts of interest to declare. Author Contributions: SMS, NN, AH, and BJR designed the study and collected the data. SMS and NN drafted the manuscript. KG provided guidance on interpreting and describing results from the stakeholder workshop and formulating recommendations. All authors interpreted results and contributed to the manuscript. Funding: This work is supported by the Agriculture and Food Research Initiative and Food and Agriculture Cyberinformatics and Tools Initiative grant #1019678 and Hatch projects #1016068 and #1015265 from the United States Department of Agriculture National Institute of Food and Agriculture. The authors would like to thank the workshop participants for sharing their viewpoints and perspectives with regard to ML models and barriers to adoption. The authors would also like to thank the editor and anonymous reviewer for their helpful comments. The authors obtained IRB approval from North Carolina State University following the stakeholder workshop in order to publish the key themes that emerged from the discussions, with no participant information disclosed. KG gratefully acknowledges the partial support of the Genetic Engineering and Society Center at NC State (<https://go.ncsu.edu/ges>).

References

- Bartling, S., Friesike, S., 2014. Opening Science: The Evolving Guide on How the Internet is Changing Research, Collaboration and Scholarly Publishing. Springer Open, Cham Heidelberg New York Dordrecht London. <https://doi.org/10.1007/978-3-319-00026-8>. Springer.
- Beck, M.W., 2018. NeuralNetTools: Visualization and analysis tools for neural networks. J. Stat. Softw. 85, 1–20. <https://doi.org/10.18637/jss.v085.i11>.
- Breiman, L., 2001. Statistical modeling: The two cultures. Stat. Sci. 16, 199–215. <https://doi.org/10.1214/ss/1009213726>.
- Donoho, D., 2017. 50 years of data science. J. Comput. Gr. Stat. 26, 745–766. <https://doi.org/10.1080/10618600.2017.1384734>.
- Faghmous, J.H., Kumar, V., 2014. A big data guide to understanding climate change: The case for theory-guided data science. Big Data 2, 155–163. <https://doi.org/10.1089/big.2014.0026>.
- Farley, S.S., Dawson, A., Goring, S.J., Williams, J.W., 2018. Situating ecology as a big-data science: Current advances, challenges, and solutions. Bioscience 68, 563–576. <https://doi.org/10.1093/biosci/biy068>.
- Faticchi, S., Vivoni, E.R., Ogden, F.L., et al., 2016. An overview of current applications, challenges, and future trends in distributed process-based models in hydrology. J. Hydrol. 537, 45–60. <https://doi.org/10.1016/j.jhydrol.2016.03.026>.
- Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., Yacalis, G., 2018. Could machine learning break the convection parameterization deadlock? Geophys. Res. Lett. 45, 5742–5751. <https://doi.org/10.1029/2018GL078202>.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. Deep Learning. MIT Press, Cambridge, MA, USA.
- Hampton, S.E., Anderson, S.S., Bagby, S.C., et al., 2015. The Tao of open science for ecology. Ecosphere 6. <https://doi.org/10.1890/es14-00402.1>.
- Hampton, S.E., Jones, M.B., Wasser, L.A., et al., 2017. Skills and knowledge for data-intensive environmental research. Bioscience 67, 546–557. <https://doi.org/10.1093/biosci/bix025>.
- Hampton, S.E., Strasser, C.A., Tewksbury, J.J., et al., 2013. Big data and the future of ecology. Front. Ecol. Environ. 11, 156–162. <https://doi.org/10.1890/120103>.
- Hanson, P.C., Stillman, A.B., Jia, X., et al., 2020. Predicting lake surface water phosphorus dynamics using process-guided machine learning. Ecol. Modell. 430, 109136. <https://doi.org/10.1016/j.ecolmodel.2020.109136>.
- Hipsey, M.R., Hamilton, D.P., Hanson, P.C., et al., 2015. Predicting the resilience and recovery of aquatic systems: A framework for model evolution within environmental observatories. Water Resour. Res. 51, 7023–7043. <https://doi.org/10.1002/2015WR017175>. Received.
- Jones, Z., Linder, F., 2015. Exploratory data analysis using random forests. In:

- Proceedings of the 73rd Annual MPSA Conference. pp. 1–31. <https://doi.org/10.21105/joss.00092>.
- Jones, Z.M., Linder, F.J., 2016. edarf: Exploratory data analysis using random forests. *J. Open Source Softw.* 1.
- Kamilaris, A., Kartakoullis, A., Prenafeta-Boldú, F.X., 2017. A review on the practice of big data analysis in agriculture. *Comput. Electron. Agric.* 143, 23–37. <https://doi.org/10.1016/j.compag.2017.09.037>.
- Karpatne, A., Atluri, G., Faghmous, J.H., et al., 2017. Theory-guided data science: A new paradigm for scientific discovery from data. *IEEE Trans. Knowl. Data Eng.* 29, 2318–2331. <https://doi.org/10.1109/TKDE.2017.2720168>.
- Lausch, A., Schmidt, A., Tischendorf, L., 2015. Data mining and linked open data - New perspectives for data analysis in environmental research. *Ecol. Modell.* 295, 5–17. <https://doi.org/10.1016/j.ecolmodel.2014.09.018>.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 521, 436–444. <https://doi.org/10.1038/nature14539>.
- Lowndes, J.S.S., Best, B.D., Scarborough, C., et al., 2017. Our path to better science in less time using open data science tools. *Nat. Ecol. Evol.* 1. <https://doi.org/10.1038/s41559-017-0160>.
- Lundberg, S.M., Erion, G., Chen, H., et al., 2020. From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* 2, 56–67. <https://doi.org/10.1038/s42256-019-0138-9>.
- NASEM, 2019. Environmental Engineering for the 21st Century: Addressing Grand Challenges. Environmental Science & Technology, Washington, DC. <https://doi.org/10.1021/acs.est.9b03244>.
- NASEM, 2018. Science Breakthroughs to Advance Food and Agricultural Research by 2030. Washington, DC. [10.17226/25059](https://doi.org/10.17226/25059).
- Norström, A.V., Cvitanovic, C., Löf, M.F., et al., 2020. Principles for knowledge co-production in sustainability research. *Nat. Sustain.* 3, pp. 182–190. <https://doi.org/10.1038/s41893-019-0448-2>.
- NRC, 2007. Models in Environmental Regulatory Decision Making. Committee on Models in the Regulatory Decision Process. National Research Council <https://doi.org/10.17226/11972>.
- Olden, J.D., Joy, M.K., Death, R.G., 2004. An accurate comparison of methods for quantifying variable importance in artificial neural networks using simulated data. *Ecol. Modell.* 178, 389–397. <https://doi.org/10.1016/j.ecolmodel.2004.03.013>.
- Olden, J.D., Lawler, J.J., Poff, N.L., 2008. Machine learning methods without tears: A primer for ecologists. *Q. Rev. Biol.* 83, 171–193. <https://doi.org/10.1086/587826>.
- Peters, D.P.C., Havstad, K.M., Cushing, J., et al., 2014. Harnessing the power of big data: infusing the scientific method with machine learning to transform ecology. *Ecosphere* 5, 67.
- Read, J.S., Jia, X., Willard, J., et al., 2019. Process-guided deep learning predictions of lake water temperature. *Water Resour. Res.* 55, 9173–9190. <https://doi.org/10.1029/2019WR024922>.
- Ribeiro, M.T., Singh, S., Guestrin, C., 2016. “Why should I trust you?” Explaining the predictions of any classifier. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. [10.1145/2939672.2939778](https://doi.org/10.1145/2939672.2939778).
- Robson, B.J., Hamilton, D.P., Webster, I.T., Chan, T., 2008. Ten steps applied to development and evaluation of process-based biogeochemical models of estuaries. *Environ. Model. Softw.* 23, 369–384. <https://doi.org/10.1016/j.envsoft.2007.05.019>.
- Rode, M., Wade, A.J., Cohen, M.J., et al., 2016. Sensors in the stream: The high-frequency wave of the present. *Environ. Sci. Technol.* 50, 10297–10307. <https://doi.org/10.1021/acs.est.6b02155>.
- Shen, C., 2018. A transdisciplinary review of deep learning research and its relevance for water resources scientists. *Water Resour. Res.* 54, 8558–8593. <https://doi.org/10.1029/2018WR022643>.
- Shen, C., Laloy, E., Elshorbagy, A., et al., 2018. HESS opinions: Incubating deep-learning-powered hydrologic science advances as a community. *Hydrol. Earth Syst. Sci.* 22, 5639–5656. <https://doi.org/10.5194/hess-22-5639-2018>.
- Thessen, A., 2016. Adoption of machine learning techniques in ecology and earth science. *One Ecosyst.* 1. <https://doi.org/10.3897/oneeco.1.e8621>.
- Wagstaff, K., 2012. Machine learning that matters, In: Proceedings of the 29th International Conference on Machine Learning. California Institute of Technology, Edinburgh, Scotland, UK. [10.1023/A:1007601113994](https://doi.org/10.1023/A:1007601113994).
- WEF, 2018. Harnessing Artificial Intelligence for the Earth, Fourth Industrial Revolution for the Earth. World Economic Forum, Geneva, Switzerland.
- Woelfle, M., Oliaro, P., Todd, M.H., 2011. Open science is a research accelerator. *Nat. Chem.* 3, 745–748. <https://doi.org/10.1038/nchem.1149>.