

stemflow: A Python Package for Adaptive Spatio-Temporal Exploratory Model

20 September 2023

Summary

Stemflow is a user-friendly python package for Adaptive Spatio-Temporal Exploratory Model (AdaSTEM, @Fink:2013) that follows the style of scikit-learn BaseEstimator class [Pedregosa:2011]. It provides one-line model creation, fitting, prediction, and evaluation. It implements spatio-temporal train-test-split and cross-validation functions. After model training, feature importance could be evaluated with spatio-temporal dynamics. Stemflow also provides functions for visualizing ensembles structured in model training and generating GIF file for predicted results to animate the spatio-temporal movement of animal population.

Statement of need

Spatio-temporal big data is an emerging but valuable resource for ecological studies as human enter the era of big data (Farley et al., 2018). Data from broad-scale survey like citizen science projects is increasingly important in modern ecological research (Dickinson et al., 2010). Intensity of survey activities grows rapidly as more people are involved in citizen science in recent years, resulted in exponential accumulation of observational data (Di Cecco et al., 2021; Sullivan et al., 2014). However, daily species observation records uploaded by non-professionals in citizen science program are known to have larger bias than professionally structured research, both in terms of data veracity and spatio-temporal balance of the datasets (Dickinson et al., 2010; Farley et al., 2018), which necessitate elaborate modeling methods to mine its merits (Dickinson et al., 2010).

Some species distribution modeling (SDM) approaches were brought forward to adjust for bias in citizen science and model on the unobserved components (Bird et al., 2014). Still, many failed to account for the autocorrelation of space and time (F. Dormann et al., 2007), which is especially crucial in modeling inherently spatio-temporal biological events with variations at different scales (Chave, 2013; Levin, 1992), such as seasonal migration. Adaptive Spatio-Temporal Exploratory Model

(AdaSTEM) is a semi-parameterized machine learning model that leverages the spatio-temporal adjacency information of sample points to model occurrence or abundance of species (Fink et al., 2013). A QuadTree algorithm (Samet, 1984) is implemented to split data into smaller spatio-temporal grids (called stixels) conditional on the data abundance, with more abundant data allowing stixels to be divided into finer resolution (up to a maximum). Stixels with data size less than a certain threshold will not be modeled; instead, these stixels will be labeled as unpredictable. This procedure controls the degree of model extrapolation and reduces overfitting. A base model is trained for each stixel, that is, targets are only modeled on their adjacent information in space and time. Splitting-training is carried out several times to generate multiple ensembles. Finally, prediction results were aggregated across these ensembles.

AdaSTEM shows the capacity of supporting large scale spatio-temporal ecological data modeling in many studies (Fink et al., 2020; Fuentes et al., 2023; La Sorte et al., 2022), especially for modeling animal abundance at different scales (Fink et al., 2013). One well-known application of AdaSTEM is the weekly abundance map of eBird Status and Trend product (Fink et al., 2022), which was widely used as data sources of abundance data of bird populations (Bird et al., 2014; Jarzyna & Stagge, 2023; Lin et al., 2022). The application of AdaSTEM could be extended to other fields with similar data structure and spatio-temporal dependence, for example, epidemiology. Despite the foreseeable significant role of spatio-temporal big data in the coming decades of scientific research, the development of tools has not necessarily kept pace.

Stemflow is positioned as a user-friendly python package to meet the need of general application of modeling spatio-temporal large datasets. Scikit-learn style object-oriented modeling pipeline enables concise model construction with compact parameterization at the user end, while the rest of the modeling procedures are carried out under the hood. Once the fitting method is called, the model class recursively splits the input training data into smaller spatio-temporal stixels using QuadTree algorithm. For each of the stixels, a base model is trained only using data falls into that stixel. Stixels are then aggregated and constitute an ensemble. In the prediction phase, stemflow queries stixels for the input data according to their spatial and temporal index, followed by corresponding base model prediction. Finally, prediction results are aggregated across ensembles to generate robust estimations (see Fink et al., 2013 and stemflow documentation for details).

For survey projects that include abundance information like eBird, the targeted modeling values are often zero-inflated, owing to the fact of low observation probability in many species. Zero-inflation could lead to poor regression model performance (Campbell, 2021). In stemflow, we implement hurdle model classes that embed two sequential models: a classifier to classify the absence and presence state, followed by a regressor to model the abundance for prediction samples classified as presence. Hurdle model classes can be conjunctively used with AdaSTEM model classes in two ways: Use hurdle model as the base model for

AdaSTEMRegressor (as in Johnston et al., 2015), or use AdaSTEMClassifier and AdaSTEMRegressor as the classifier and regressor in hurdle model. We demonstrate the comparison of these two architectures in stemflow documentation.

One advantage of applying stemflow in scikit-learn style is that there is a variety of “base models” to choose from scikit-learn or scikit-learn-style repertoire. The choices vary from linear models to boosting and bagging tree-based models. Maxent model (C. B. Anderson, 2023) is also supported to play the role of “base model”, which largely expands the potential application for presence-only modeling (see documentation).

While there exists mounting open source packages for species distribution modeling (mostly in R, Norberg et al., 2019; and one in Python, C. B. Anderson, 2023), most of them solely leverage environmental variables and do not support integration of spatio-temporal information during model construction (but see C. B. Anderson, 2023; S. C. Anderson et al., 2022; Dobson et al., 2023). This disadvantage is commonly criticized along with the overconfidence of the model extrapolation capacity both for Maxent-based and ensemble-based models (A. Lee-Yaw et al., 2022). To our knowledge, stemflow is the first package designed to solve the spatio-temporal dependency conjugated with bias in data abundance distribution in species distribution modeling. With the rapid accumulation of data and development of machine learning techniques, stemflow will show potentials in spatio-temporal modeling, and could be applied to other fields (e.g., epidemiology and weather prediction) in future.

Acknowledgements

This project was based upon work supported by National Natural Science Foundation of China grants 32125005 (to X.Z.), 32270455 (to Z.G.), CAS Project for Young Scientists in Basic Research (YSBR-097 to X.Z. and Z.G.). We thank Dr. Daniel Fink at Cornell Lab of Ornithology for the R scripts of STEM model and QuadTree which inspired our implementation in python.

References