

stacks: Stacked Ensemble Modeling with Tidy Data Principles

Simon P. Couch¹ and Max Kuhn¹

¹ RStudio PBC

DOI:

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Submitted:

Published:

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

Summary

Model stacking is an ensemble modeling technique that involves training a model to combine the outputs of many constituent statistical models. `{stacks}` is a free and open-source R software package for stacked ensemble modeling that is consistent with tidy data principles. The package’s functionality is closely aligned with the `{tidymodels}`, a collection of packages providing a unified interface to a diverse set of statistical modeling techniques. Beyond simply providing a mathematically robust interface to build stacked ensemble models, `{stacks}` adheres to a consistent grammar in order to interface with two object classes that promote an intuitive understanding of the underlying implementation.

Statement of Need

Model stacking, one approach to creating an *ensemble model*, is an increasingly popular approach in statistical and machine learning that involves training a model to generate predictions informed by many constituent models (“members”). Model stacking has been shown to increase predictive performance in a variety of settings and has thus greatly risen in popularity in recent years. In a spring 2020 community survey, the advanced statistical modeling technique was ranked as the first priority for future development among users of the `{tidymodels}`, a burgeoning software ecosystem for tidyverse-aligned predictive and inferential modeling used across many modern research and industrial applications (M. Kuhn & Wickham, 2022). `{stacks}` introduces model stacking to the `{tidymodels}`.

Packages implementing methods for predictive and inferential modeling in R are highly variable in their interfaces. The structure of inputted data, argument names, expected argument types, argument orders, output types, and spelling cases varies widely both within and among packages. This diversity in approaches obscures the intuition shared among common modeling procedures, makes details of usage difficult to remember, and prevents an expressive and idiomatic coding style. In contrast, `{stacks}` utilizes the consistent and unified interface of the `{tidymodels}` packages to implement a generalized and concise grammar for model stacking. The package supports ensembling using any member model type, cross-validation scheme, and error metric implemented in—or in alignment with—the `{tidymodels}` (M. Kuhn & Wickham, 2022).

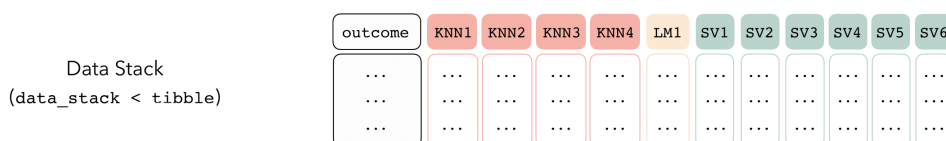
The principled and generalized approach of the package lends itself to diverse applications of predictive modeling. This capability, for one, was recognized with the receipt of the American Statistical Association’s 2021 John M. Chambers Award, awarded “for the development and implementation of computational tools for the statistical profession by a graduate or undergraduate student” (ASA Section on Statistical Computing, 2021). The package’s functionality has also been shared in venues such as R/Pharma 2020,

rstudio::global(2021), and the 2021 Joint Statistical Meetings. To date, the package has been downloaded more than 20,000 times, evidencing its key contribution to a software ecosystem utilized in diverse research contexts.

Underlying Principles

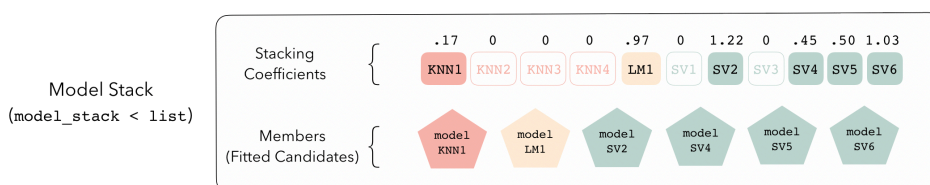
At a high level, stacked ensemble models in `{stacks}` are formed from *model definitions*, which specify a set of instructions to fit a model or set of models using `{parsnip}`, `{recipes}`, and `{workflows}` (Max Kuhn & Vaughan, 2022; Max Kuhn & Wickham, 2022; Vaughan, 2022). Each model specified in a model definition is referred to as a *candidate member* in the package—model definitions can result in multiple candidate members when several possible hyperparameter values are being optimized over (e.g. the number of neighbors k or the precision parameter for the support vector machine’s radial basis function).

To be used in the same model stack, all model definitions must share a *resampling object*, as defined in the `{rsample}` package, such as a cross-fold validation or set of bootstrap samples (Silge, Chow, Kuhn, & Wickham, 2022). After initializing a `data_stack` object with the `stacks()` function, model definitions can be iteratively added to the data stack with `add_candidates()`. This function collates the predictions on the assessment set, a subset of the training data used for preliminary model validation, from each candidate specified in the model definitions to the data stack.



The above diagram represents a data stack containing 11 candidate members, where 4 come from a shared K-nearest neighbors model definition, 1 arises from a linear model, and 6 come from a shared support vector machine model definition.

After all candidates are collated to the data stack, the user can fit the meta-learner using the `blend_predictions()` function, which fits an elastic net model on the data stack, predicting the true outcome using the predictions from each of the candidate members (Hastie, Tibshirani, & Wainwright, 2015). The coefficients of this model form the *stacking coefficients*, which are the weightings for each of the member models in ultimately predicting the outcome. Candidate members with non-zero stacking coefficients are *members*, and must be fitted on the entire training set with the `fit_members()` function. This function outputs a `model_stack` object, and is ready to predict on new data.



In addition to the aforementioned core verbs, the package supplies several helper functions to interface more effectively with model stacks. Notably, `collect_parameters()` juxtaposes model parameters with their stacking coefficients, and an `autoplot()` S3 method provides model diagnostic visualizations using `{ggplot2}` (Wickham, 2016).

Package documentation, including long-form vignettes and relevant links, is available on the package website stacks.tidymodels.org.

Comparison to Other Software

A number of software packages on the Comprehensive R Archive Network share functionality with `{stacks}` (R Core Team, 2022). Notably, of course, the package integrates tightly with other packages in the `{tidymodels}` ecosystem, such as `{tune}`, `{parsnip}`, `{rsample}`, and `{recipes}` (Max Kuhn, 2022; Max Kuhn & Vaughan, 2022; M. Kuhn & Wickham, 2022; Max Kuhn & Wickham, 2022; Silge et al., 2022).

There are also other packages providing implementations of model ensembling in R. The `{h2o}` R package ports functionality from the H2O modeling ecosystem to R via a REST API, including an implementation of ensembling. Notably, the `{h2o}` R package supports a generalized meta-learner, rather than only elastic nets and their special cases like `{stacks}`, and uses a more abbreviated grammar as a result (LeDell et al., 2022; LeDell & Poirier, 2020). Similarly, the `{mlr3pipelines}` package implements an interface for ensembling models defined in the `{mlr3}` modeling ecosystem with similar interface differences as those noted for `{h2o}` (Binder et al., 2021; Lang et al., 2019). The `{SuperLearner}` package also provides an implementation of model ensembling. The package provides its own wrappers for member model types and error metrics that provide less granular control than `{stacks}`' dependencies, though also supplies a wider array of meta-learner options than `{stacks}`, again resulting in a more abbreviated interface (Laan, Polley, & Hubbard, 2007). Beyond the differences noted, each of these packages differ in the classes of model types, error metrics, cross-validation schemes, and meta-learners supported, as well as the modeling behaviors encouraged in their interfaces.

Beyond the R language, a number of modeling ecosystems implement interfaces to ensembling. In Python, notably, the Keras and Scikit-learn libraries supply methods for model ensembling (Chollet et al., 2015; Pedregosa et al., 2011). Many modeling ecosystems utilized in multi-language ports include model ensembling in implementations of automated machine learning—notable examples include the Azure Machine Learning platform and (previously mentioned) H2O AutoML (Barnes, 2015; LeDell & Poirier, 2020). Model ensembling, generally, is an actively researched topic with diverse applications and implementations—see Dong et al. for a review of contemporary ensembling research and references to software implementations (Dong, Yu, Cao, Shi, & Ma, 2020). Again, each of these ecosystems vary greatly in their degree of extensibility and the modeling behaviors that their interfaces encourage and accommodate.

Acknowledgements

We acknowledge contributions from Julia Silge, Hannah Frick, and Asmae Toumi. Development of this package was supported by RStudio PBC, the Reed College Mathematics Department, and the American Statistical Association Sections on Statistical Computing & Statistical Graphics.

References

ASA Section on Statistical Computing. (2021). *John M. Chambers Statistical Software Award*. Retrieved from <https://community.amstat.org/jointscsg-section/awards/john-m-chambers>

- Barnes, J. (2015). Azure Machine Learning. *Microsoft Azure Essentials*. ISBN: 978-0-7356-9817-8; Microsoft Press.
- Binder, M., Pfisterer, F., Lang, M., Schneider, L., Kotthoff, L., & Bischl, B. (2021). mlr3pipelines - flexible machine learning pipelines in R. *Journal of Machine Learning Research*, 22(184), 1–7. Retrieved from <https://jmlr.org/papers/v22/21-0281.html>
- Chollet, F. et al. (2015). Keras. <https://keras.io>.
- Dong, X., Yu, Z., Cao, W., Shi, Y., & Ma, Q. (2020). A survey on ensemble learning. *Frontiers of Computer Science*, 14(2), 241–258.
- Hastie, T., Tibshirani, R., & Wainwright, M. (2015). Statistical learning with sparsity. *Monographs on statistics and applied probability*, 143, 143.
- Kuhn, Max. (2022). *Tune: Tidy tuning tools*. Retrieved from <https://CRAN.R-project.org/package=tune>
- Kuhn, Max, & Vaughan, D. (2022). *Parsnip: A common API to modeling and analysis functions*. Retrieved from <https://CRAN.R-project.org/package=parsnip>
- Kuhn, M., & Wickham, H. (2022). tidymodels: A collection of packages for modeling and machine learning using tidyverse principles. *Boston, MA, USA*. Retrieved from <https://tidymodels.org>
- Kuhn, Max, & Wickham, H. (2022). *Recipes: Preprocessing and feature engineering steps for modeling*. Retrieved from <https://CRAN.R-project.org/package=recipes>
- Laan, M. J. van der, Polley, E. C., & Hubbard, A. E. (2007). Super Learner. *Statistical Applications in Genetics and Molecular Biology*, 6(1). doi:10.2202/1544-6115.1309
- Lang, M., Binder, M., Richter, J., Schratz, P., Pfisterer, F., Coors, S., Au, Q., et al. (2019). mlr3: A modern object-oriented machine learning framework in R. *Journal of Open Source Software*. doi:10.21105/joss.01903
- LeDell, E., Gill, N., Aiello, S., Fu, A., Candel, A., Click, C., Kraljevic, T., et al. (2022). *h2o: R interface for the 'H2O' scalable machine learning platform*. Retrieved from <https://CRAN.R-project.org/package=h2o>
- LeDell, E., & Poirier, S. (2020). h2o AutoML: Scalable automatic machine learning. *Proceedings of the AutoML workshop at ICML* (Vol. 2020).
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., et al. (2011). Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.*, 12(null), 2825–2830.
- R Core Team. (2022). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from <https://www.R-project.org/>
- Silge, J., Chow, F., Kuhn, M., & Wickham, H. (2022). *Rsample: General resampling infrastructure*. Retrieved from <https://CRAN.R-project.org/package=rsample>
- Vaughan, D. (2022). *Workflows: Modeling workflows*. Retrieved from <https://CRAN.R-project.org/package=workflows>
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. Retrieved from <https://ggplot2.tidyverse.org>