

Assessing Feature Importance in Non-Linear Regression: An Ecology Perspective

DEPARTMENT OF BIOLOGY

W

Anthony F. Cannistra
tonycan@uw.edu

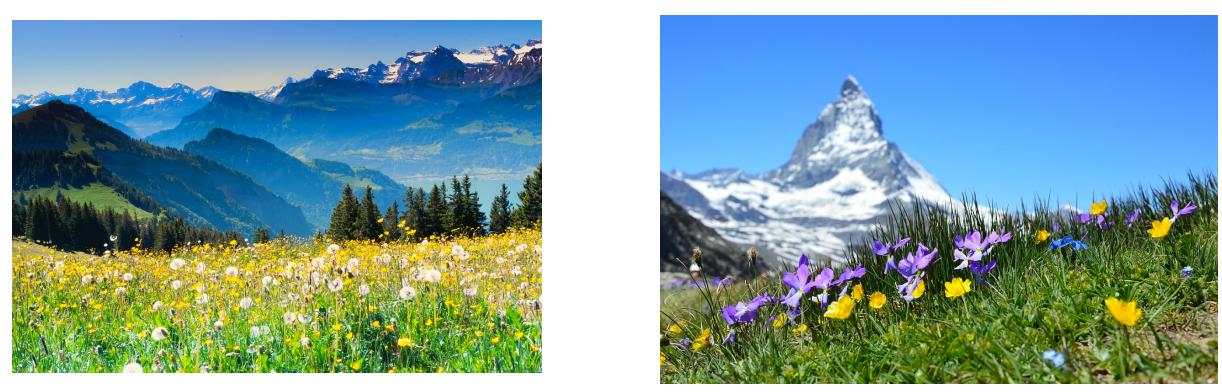
Motivation:

- For years, linear regression has been used by ecologists to study and quantify observed ecological relationships.
- These tools have been useful in understanding ecological mechanism, but our societal mandate to **predict future ecological state** requires models which offer more predictive power than a simple linear method can (there's reason to believe many relationships are nonlinear).
- To maintain scientific and practical credibility, these models must be at least to some extent *interpretable* to both researchers and the users of their predictions.
- Simple methods to capture the nonlinearities of many ecological relationships with provable properties are kernel regression and SVMs.

Key Question: What methods exist to extract variable importance from a non-linear regression, and how do the results therein compare to coefficients from a standard linear regression model? **Do the non-linear methods still capture theoretically-rigorous biology?**

Data:

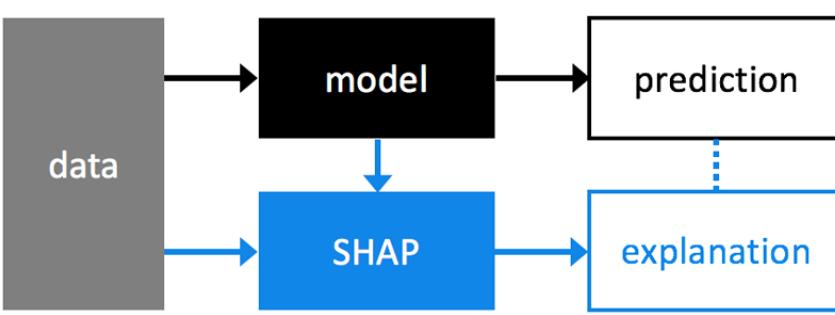
- Angert et al. 2011** gathered **functional trait data** for Swiss alpine plants ($N = 133$, **Holzinger et al. 2008**) and computed the *observed elevational range shift* from historical records (i.e.: how much has the maximum elevation at which these plants are observed shifted in the past century?). There are $d = 32$ observed traits.
- This is one of two datasets analyzed.



Methods:

- We compare **Ordinary Least Squares' coefficient values** to:
 - Kernel Regression (RBF, Polynomial) coefficients
 - Ridge Regression coefficients
 - Random Forest Gini Scores
 - Shapley Values¹** for SVM (RBF, Polynomial)

1: The Shapley additive feature importance measure, defined in **Lundberg and Lee 2017**, from cooperative game theory.

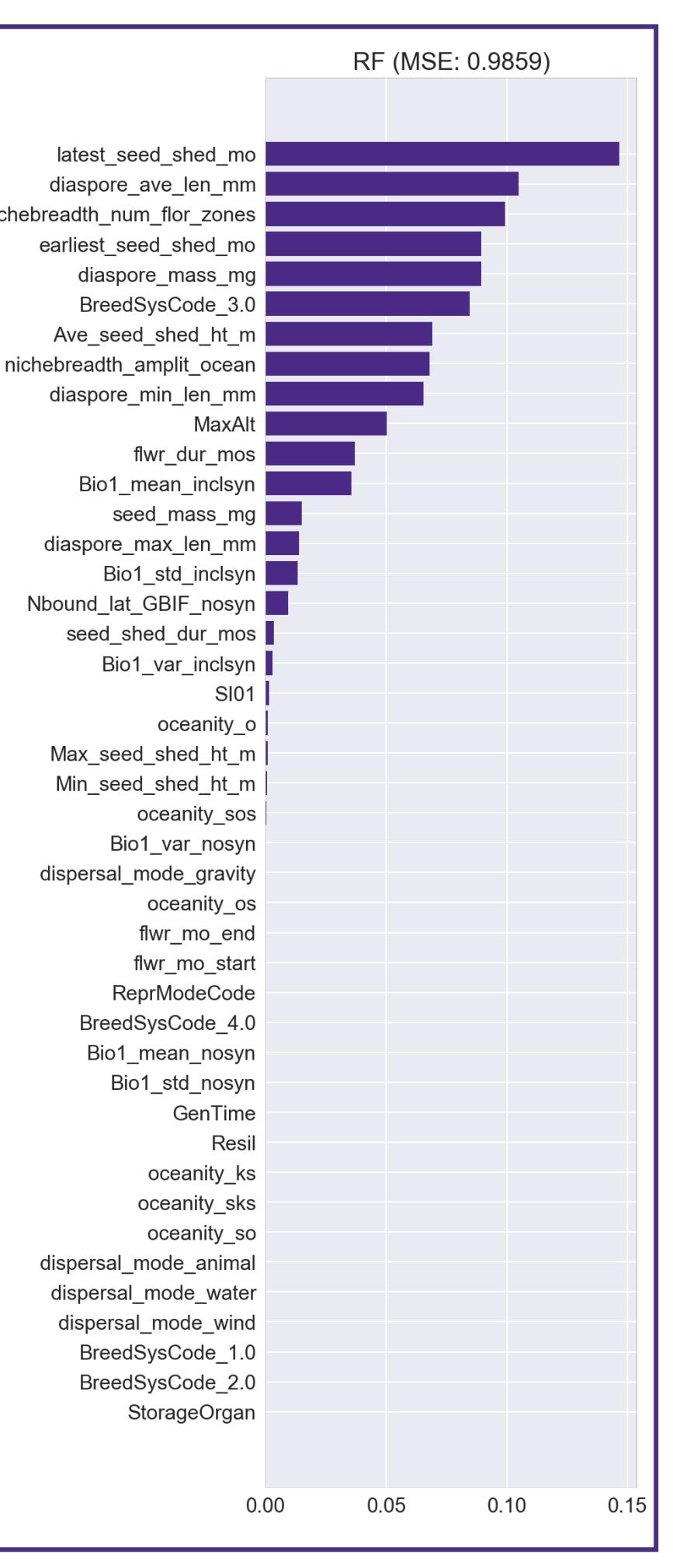
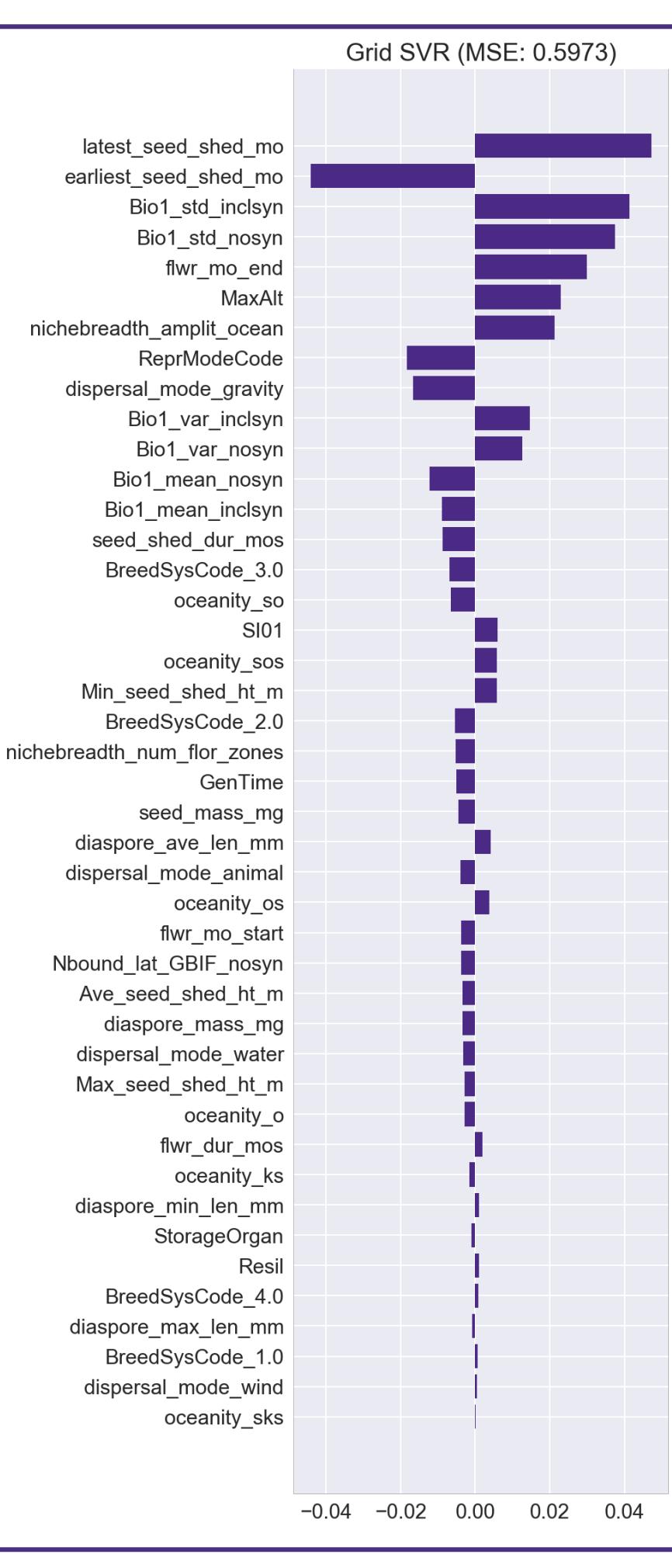
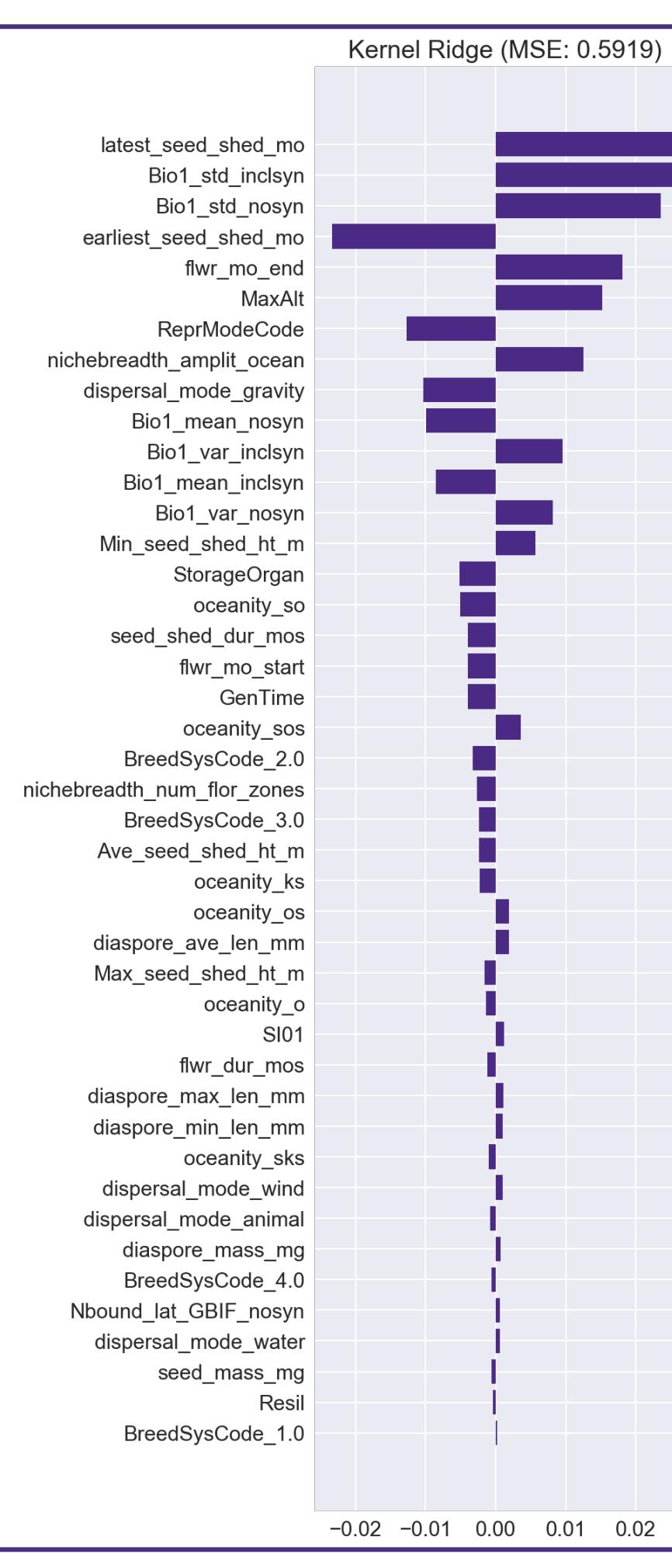
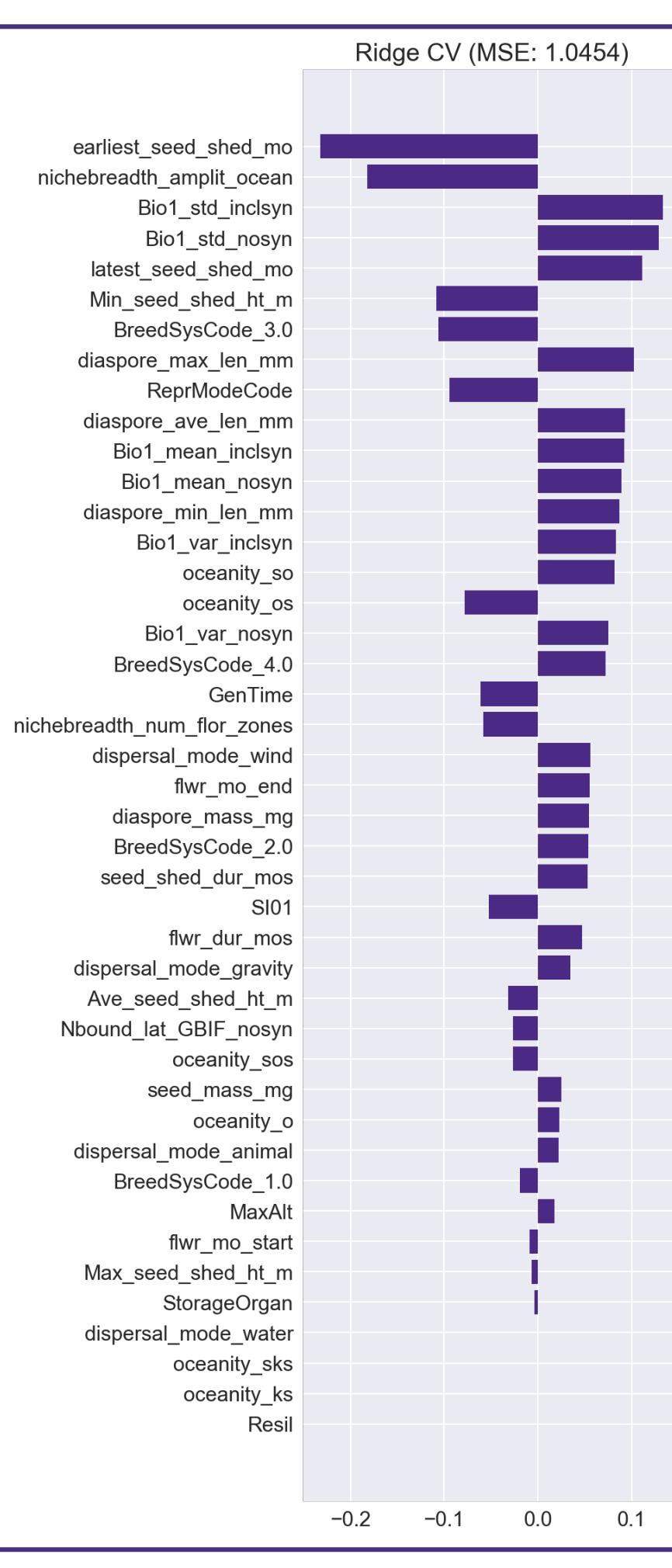
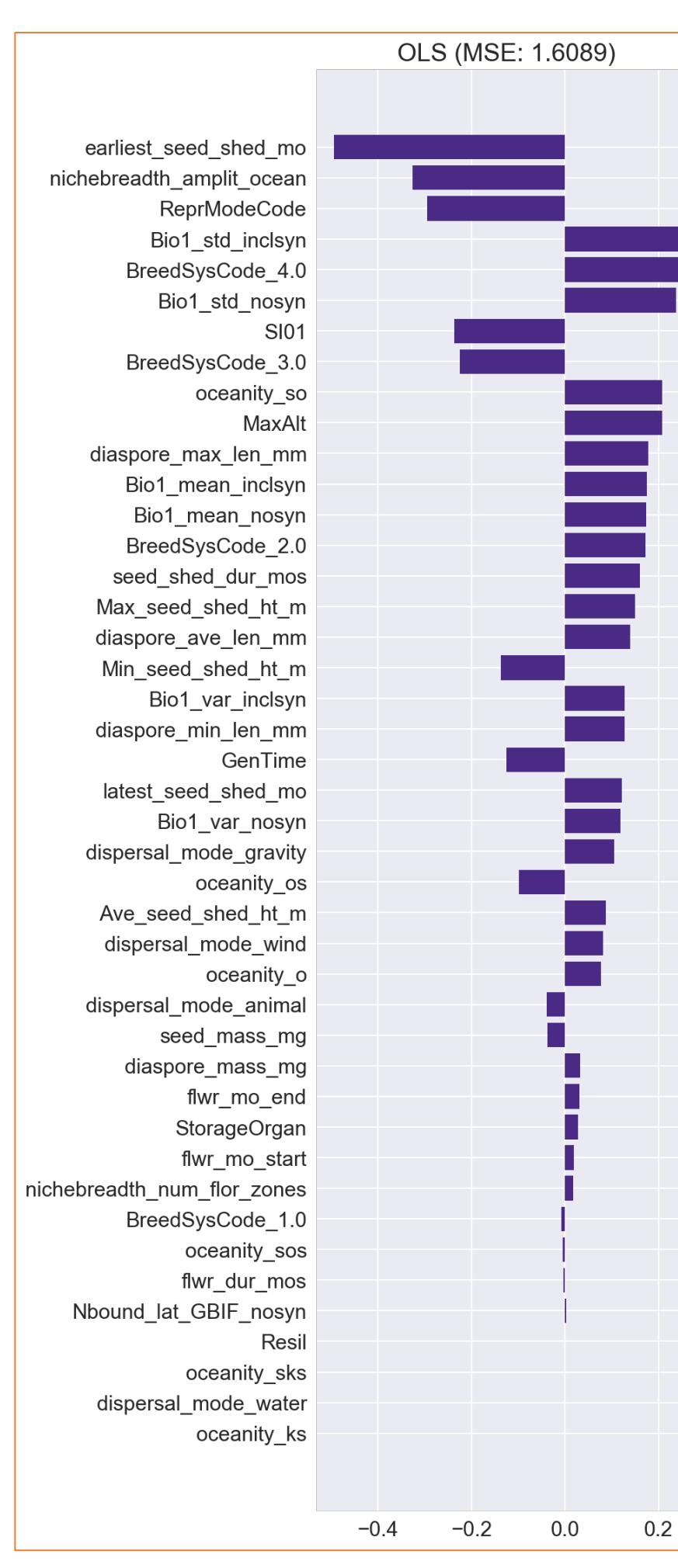


Implementation:

- We remove samples (i.e., plant species) which are missing any feature data, which brings us to $N = 20$.
- We one-hot encode categorical features, $d = 43$.
- We center the numeric features with 0 mean and unit variance for regression.

References:

- Angert, A. et al., 2011. "Do species' traits predict recent shifts at expanding range edges?" *Ecology Letters*, **14**: 677–689.
Holzinger, B. et al., 2008. "Changes in plant species richness over the last century in the eastern Swiss Alps: elevational gradient, bedrock effects and migration rates." *Plant Ecol.*, **195**, 179–196.
Lundberg, S. and Lee, S., 2017. "A unified approach to interpreting model predictions" *NIPS 2017*. arXiv:1705.07874



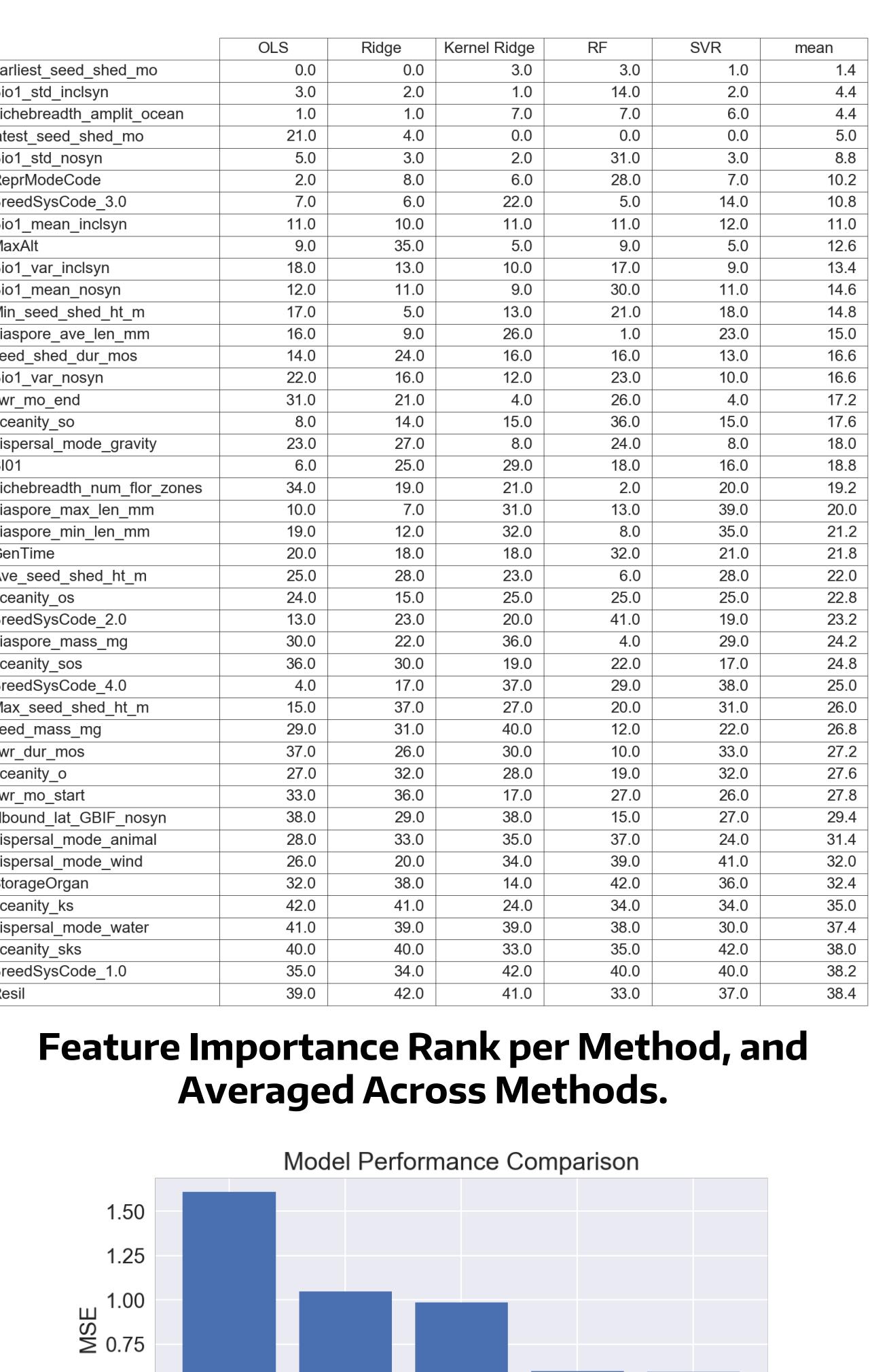
OLS Coefficients

Ridge Regression Coefs.

Kernel Regression SHAP

SVM SHAP

Random Forest Gini Scores



Feature Importance Rank per Method, and Averaged Across Methods.



Results | Discussion | Conclusions:

- Across these methods, **dispersal-related traits** like {`latest_seed_shed_mo`, `earliest_seed_shed_mo`} consistently rank highest in importance across these methods, as measured by the absolute value of the coefficients. For all methods except for the Ridge and OLS, `latest_seed_shed_mo` was the variable with the highest importance.
- These results are consistent with Angert et al.'s 2011 analysis, wherein seed dispersal period duration was the most important predictor variable.
- This result indicates that the predictive performance gains (see MSE figure above) from the SVM/Kernel Ridge regression maintain their theoretical biological underpinnings while better representing nonlinearities in the relationship between species range shift and trait values.