

# STAR methods tables

<b>Methods S1. Definition of terms, related to Table 1, Figure 1(b), and Figure 1(c)</b>	
<b>Richness</b>	The count of species or KOs detected in a sample
<b>Mean absolute error</b>	$\frac{\sum_{i=1}^n  y_i - x_i }{n}$
$R^2$	<p>The proportion of the variance in the dependent variable that is predictable from the independent variables.</p> $1 - \frac{RSS}{TSS}$ <p>Where RSS is the sum of squares of residuals, and TSS is the total sum of squares.</p>
<b>Shannon's diversity</b>	$-\sum_{i=1}^R p_i \ln(p_i)$ <p>Where <math>p_i</math> is the empirical probability of detecting species <math>i</math> in the dataset.</p>
<b>Pielou's Evenness index</b>	$\frac{H}{H_{Max}}$ <p>Where H is shannon diversity and <math>H_{MAX}</math> is the maximum possible achievable shannon entropy.</p>

<b>Methods S2 related to Table 1, Figure 1(b), and Figure 1(c)</b>		
<b>Library</b>	<b>Relevant functions/modules</b>	<b>Description</b>
<a href="https://scikit-learn.org/">https://scikit-learn.org/</a> Version 0.22.0	<ol style="list-style-type: none"> <li>1. sklearn.linear_model.ElasticNet</li> <li>2. sklearn.model_selection.GridSearchCV</li> <li>3. sklearn.preprocessing.PowerTransformer</li> <li>4. sklearn.pipeline.Pipeline</li> <li>5. sklearn.metrics.r2_score</li> <li>6. sklearn.metrics.mean_absolute_error</li> <li>7. sklearn.metrics.train_test_split</li> <li>8. sklearn.model_selection.RepeatedStratifiedKfold</li> <li>9. sklearn.ensemble.RandomForestRegressor</li> </ol>	<ol style="list-style-type: none"> <li>1. Elastic net implementation</li> <li>2. Used for hyperparameter optimization</li> <li>3. Power transformation implementation</li> <li>4. Utility to chain commands that need to be applied during training</li> <li>5. Coefficient of determination</li> <li>6. MAE</li> <li>7. Used to create train/test/validation data</li> <li>8. Used for nested cross validation</li> <li>9. Random forest implementation</li> </ol>
<a href="http://scikit-bio.org/">http://scikit-bio.org/</a> Version 0.5.6	<ol style="list-style-type: none"> <li>1. skbio.stats.composition.multiplicative_replacement</li> <li>2. skbio.stats.composition.clr</li> </ol>	<ol style="list-style-type: none"> <li>1. Implementation of imputation using multiplicative replacement</li> <li>2. CLR transformation implementation</li> </ol>
Xgboost Version 1.0.2	<ol style="list-style-type: none"> <li>1. XGBRegressor</li> </ol>	<ol style="list-style-type: none"> <li>1. Gradient boosting machine implementation</li> </ol>
PyTorch Version 1.1.0	<ol style="list-style-type: none"> <li>1. torch.utils.data: DataLoader, TensorDataset</li> <li>2. torch.autograd.Variable</li> <li>3. 3. torch.optim.*</li> <li>4. 4. torch.nn.functional</li> <li>5. 5. torch.nn</li> </ol>	<ol style="list-style-type: none"> <li>1. Dataset utilities</li> <li>2. Autodiff utilities</li> <li>3. Optimization</li> <li>4. Functional neural net API</li> <li>5. Neural network functions</li> </ol>

<b>Methods S3 related to Table 1, Figure 1(b), and Figure 1(c)</b>		
<b>Model class</b>	<b>Hyperparameter Space</b>	<b>Model selection method</b>
Neural net	1 and 2 hidden layer NNs varying hidden dimension from 50 to 1500.	Grid search using held out validation set.
Random forest	n_estimators = [50, 100, 200, 300, 500] max_depth = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, None] min_samples_split = [2, 5, 10] min_samples_leaf = [1, 2, 4]	Random search + CV
Gradient machine boosting	learning_rate = [0.05, 0.10, 0.15, 0.20, 0.25, 0.30] , max_depth = [ 3, 4, 5, 6, 8, 10, 12, 15] min_child_weight=[ 1, 3,5,7], gamma = [ 0.0, 0.1, 0.2, 0.3, 0.4] colsample_bytree = [0.3, 0.4, 0.5, 0.7 ] n_estimators =[50, 100, 200, 300, 500, 600]	Random search +CV

<b>Methods S4 related to Table 1, Figure 1(b), and Figure 1(c)</b>		
<b>Model class</b>	<b>Feature type</b>	<b>Selected model</b>
Neural net	Blood	2 hidden layers, hidden dim =1000
	Stool	1 hidden layer, hidden dim=500
Random forest	Blood	n_estimators = 200 max_depth = None min_samples_split = 10 min_samples_leaf = 2
	Stool	n_estimators = 100 max_depth = 5 min_samples_split = 2 min_samples_leaf = 1
Gradient boosting machine	Blood	max_depth = 5 min_child_weight= 1 gamma =0.3 colsample_bytree = 0.5 n_estimators = 500
	Stool	max_depth = 4 min_child_weight= 5 gamma = 0.0 colsample_bytree = 0.3 n_estimators = 500