



modelimportance: An R package for evaluating model importance within an ensemble

Minsu Kim 

University of Massachusetts Amherst

Nicholas G. Reich 

University of Massachusetts Amherst

Abstract

Ensemble forecasts are commonly used to support decision-making and policy planning across various fields because they often offer improved accuracy and stability compared to individual models. As each model has its own unique characteristics, understanding and measuring the value each constituent model adds to the overall accuracy of the ensemble is of great interest to building effective ensembles. The R package **modelimportance** provides tools to quantify how each component model contributes to the accuracy of ensemble performance for different forecast types. It supports multiple functionalities; it allows users to specify which ensemble approach to implement and which model importance metric to use. Additionally, the software offers customizable options for handling missing values. These features enable the package to serve as a versatile tool for researchers and practitioners aiming to construct an effective ensemble model across a wide range of forecasting tasks.

Keywords: ensemble, forecast, prediction, model importance, Shapley value, R.

1. Introduction

Ensemble forecasting is a method to produce a single, consolidated prediction by combining forecasts generated from different models. This technique leverages the strengths of individual models while mitigating their weaknesses, leading to more robust and accurate predictions (Gneiting and Raftery 2005; Hastie *et al.* 2001; Lutz *et al.* 2019; Viboud *et al.* 2018). Specifically, ensembles aggregate diverse insights from various models and handle variance and bias more effectively by averaging out the biases and variations of individual models, which can reduce prediction errors and improve overall model performance. Due to these enhanced prediction accuracy and robustness, it is widely used across various domains such as weather

forecasting (Jordan A. Guerra *et al.* 2020; Gneiting and Raftery 2005), financial modeling (Sun *et al.* 2020; He *et al.* 2023), and infectious disease outbreak forecasting (Ray and Reich 2018; Reich *et al.* 2019) to improve decision-making and achieve more reliable predictions. For example, throughout the COVID-19 pandemic, the US COVID-19 Forecast Hub collected individual models developed by over 90 different research groups and built a probabilistic ensemble forecasting model for COVID-19 cases, hospitalizations, and deaths in the US based on those models' predictions, which served as the official short-term forecasts for the US Centers for Disease Control and Prevention (CDC) (Kim *et al.* 2024).

The quality of forecasts is assessed by evaluating their error, bias, sharpness, and/or calibration using different scoring metrics. The selection of the scoring metrics depends on the type of forecast, such as point forecasts and probabilistic forecasts, and their corresponding formats, such as median/mean, quantiles, and predictive cumulative distribution functions. The mean absolute error (MAE) and the mean squared error (MSE) are the commonly used measures for point forecast accuracy. They provide a direct assessment of the accuracy of specific forecasted values by calculating the average magnitude of errors. Probabilistic forecasts are typically assessed using proper scoring rules, which consider the uncertainty and variability in predictions, providing concise evaluations through numerical scores (Gneiting and Raftery 2007). Some examples include the weighted interval score (WIS) for the quantile-based forecasts and the continuous ranked probability score (CRPS) for the forecasts taking the form of predictive cumulative distribution functions (Bracher *et al.* 2021).

Several R packages have been developed for this purpose. To name a few, the **forecast** package (Hyndman and Khandakar 2008) is widely used for univariate time series forecasting and includes functions for accuracy measurement. The **Metrics** (Hamner and Frasco 2018) and **MLmetrics** (Yan 2024) provide a wide range of performance metrics specifically designed for evaluating machine learning models. The **scoringRules** (Jordan *et al.* 2019) package offers a comprehensive set of proper scoring rules for evaluating probabilistic forecasts and supports both univariate and multivariate settings. The **scoringutils** (Bosse *et al.* 2022) package adds to the functionality provided by **scoringRules**, offering additional features that make it more useful for certain tasks, such as summarizing, comparing, and visualizing forecast performance. Additionally, the **scoringutils** supports forecasts represented by predictive samples or quantiles of predictive distributions, allowing for the assessment of any forecast type, even when a closed-form expression for a parametric distribution is not available. These packages have been valuable to evaluate individual models as independent entities, using performance metrics selected for each specific situation or problem type. However, they do not measure the individual models' contributions to the enhanced predictive accuracy when used as part of an ensemble. Kim *et al.* (Kim *et al.* 2024) demonstrate that a model's individual performance does not necessarily correspond to its contribution as a component within an ensemble. Our developed package introduces this capability. The **modelimportance** package provides tools to evaluate the role of each model as an ensemble member within an ensemble model, rather than focusing on the individual predictive performance per se.

In ensemble forecasting, certain models contribute more significantly to the overall predictions than others. Assessing the impact of each component model on ensemble predictions is methodologically similar to determining variable importance in traditional regression and machine learning models, where variable importance measures evaluate how much individual variables enhance the model's predictive performance. R packages that implement these func-

tions include **randomForest** (Liaw and Wiener 2002), **caret** (Kuhn 2008), **xgboost** (Chen *et al.* 2024), and **gbm** (Ridgeway and Developers 2024), each providing variable importance measures for different types of models: random forest models, general machine learning models, extreme gradient boosting models, and generalized boosted regression models, respectively. Similarly, the tools in the **modelimportance** package quantify the contribution of each component model within an ensemble to the ensemble model’s predictive performance. They assign numerical scores to each model based on a selected metric that measures forecast accuracy, depending on the forecast type.

These capabilities provide unique support for hub organizers, such as the CDC in the US and the European Centre for Disease Prevention and Control in the EU. By enabling precise evaluation of each model’s contribution, the **modelimportance** package helps these organizations gather and utilize forecasts from various models to create more effective ensemble forecasts. This, in turn, facilitates better communication with the public and decision-makers and enhances the overall decision-making process. Specifically, **modelimportance** is incorporated into the ‘hubverse’, which is a collection of open-source software and data tools developed to promote collaborative modeling hub efforts and reduce the effort required to set up and operate them (Consortium of Infectious Disease Modeling Hubs 2024). This package adheres to the model output formats specified by the hubverse convention, which enables seamless integration and interoperability with other forecasting tools and systems.

The paper proceeds as follows. In Section 2, we address the model output formats proposed within the hubverse framework and the structure of forecasts, followed by a motivating example. Section 3 presents two algorithms implemented in **modelimportance** for calculating the model importance metric: leave-one-model-out and leave-all-subsets-of-models-out. We demonstrate the various functionalities **modelimportance** supports in Section 4 and give some examples in Section 5. We close this paper with some concluding remarks and a discussion of possible extensions.

2. Data

2.1. Model output format

Model outputs are structured in a tabular format designed specifically for predictions. In the hubverse standard, each row represents an individual prediction for a single task, and its details are described in multiple columns through which one can identify the model IDs, task characteristics, prediction representation type, and predicted values (Shandross *et al.* 2024). In Table 1, for example, the `model_id` column contains the uniquely identified name of the model that produced the prediction in each row. Task characteristics are represented by `reference_date`, `target`, `horizon`, `location`, and `target_end_date` columns, collectively referred to as the task ID columns. The prediction representation type is specified in the `output_type` and `output_type_id` columns. This example illustrates short-term forecasts of incident influenza hospitalizations in the US for Massachusetts (FIPS code 25), generated by the model ‘Flusight-baseline’ on November 19, 2022. The forecasts are provided in seven quantiles (0.05, 0.1, 0.25, 0.5, 0.75, 0.9, 0.95) for each target end date.

model_id	reference_date	target	horizon	location	target_end_date	output_type	output_type_id	value
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.05	22
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.1	31
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.25	45
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.5	51
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.75	57
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.9	71
Flusight-baseline	2022-11-19	wk inc flu hosp	0	25	2022-11-19	quantile	0.95	80
Flusight-baseline	2022-11-19	wk inc flu hosp	1	25	2022-11-26	quantile	0.05	5
Flusight-baseline	2022-11-19	wk inc flu hosp	1	25	2022-11-26	quantile	0.1	21
Flusight-baseline	2022-11-19	wk inc flu hosp	1	25	2022-11-26	quantile	0.25	38

Table 1: Example of the model output for incident influenza hospitalizations extracted from `forecast_outputs` data in the **hubExamples** package.

2.2. Structure of forecasts

The output of the forecasting model can be represented in various types. Generally, quantitative forecasts can be categorized into either point forecasts or probabilistic forecasts. For a specific task, point forecasts, represented by a single predicted value, provide a clear and concise prediction, making them easy to interpret and communicate. Probabilistic forecasts, on the other hand, describe the likelihood of different outcomes, conveying the uncertainty inherent in the prediction. These forecasts are represented by a probability distribution over possible future values in various ways, such as probability mass functions (pmf), cumulative distribution functions (cdf), or probability quantiles (or intervals).

The `output_type` and `output_type_id` columns in the model output format, as defined by the hubverse convention, specify the forecast structure. For example, `output_type` may be ‘mean’ or ‘median’ for point forecasts, and ‘pmf’, ‘cdf’, or ‘quantile’ for probabilistic forecasts. The `output_type_id` column provides additional details for probabilistic forecasts; for instance, when the `output_type` is ‘quantile’, the `output_type_id` identifies specific quantile levels, as illustrated in Table 1. Different output types correspond to different scoring rules for evaluating a model’s prediction performance. Table 2 presents the output types and their associated scoring rules supported by the **modelimportance** package.

Output Type	Scoring Rule	Description
mean	MSE	Evaluate using the mean squared error (MSE)
median	MAE	Evaluate using the mean absolute error (MAE)
quantile	WIS	Evaluate using the weighted interval score (WIS)
pmf	Log Score	Evaluate using the logarithm of the probability assigned to the true outcome (LogScore)

Table 2: Pairs of output types and their associated scoring rules for evaluating prediction performance.

3. Algorithms

Algorithm 1 Importance score using leave one model out algorithm

Input:

- Set of n individual models, $A = \{1, 2, \dots, n\}$, and their forecasts $\{F^1, F^2, \dots, F^n\}$
- Truth value, y
- Function g to build an ensemble
- Scoring rule ψ to evaluate forecast skills

Output: importance metric of model i , $\phi^{i,\text{lomo}}$

- 1: Create an ensemble forecast F^A using g : $F^A \leftarrow g(\{F^1, F^2, \dots, F^n\})$
 - 2: Evaluate F^A using ψ : $\psi(F^A, y)$.
 - 3: **for** each i , ($i = 1, 2, \dots, n$) **do**
 - 4: $F^{A^{-i}} \leftarrow g(\{F^j | j \neq i, j \in A\})$
 - 5: Compute $\psi(F^{A^{-i}}, y)$.
 - 6: $\phi^{i,\text{lomo}} \leftarrow \psi(F^{A^{-i}}, y) - \psi(F^A, y)$
 - 7: **end for**
-

This section provides a brief description of the leave one model out (LOMO) and leave all subsets of models out (LASOMO) algorithms. (Details can be found in (Kim *et al.* 2024))

LOMO involves creating an ensemble by excluding one component model from the entire set of models. Let A be a set of n models and F^i be a forecast produced by model i , where $i = 1, 2, \dots, n$. Each ensemble excludes exactly one model while including all the others. Denoting by $F^{A^{-i}}$ an ensemble forecast constructed without F^i and by F^A the ensemble forecast built from the entire set of models, the importance score using LOMO is calculated by the difference of measures of these two ensemble performances, $F^{A^{-i}}$ and F^A . For example, when evaluating model 1 within an ensemble of three models ($n = 3$), LOMO creates an ensemble forecast $F^{\{2,3\}}$ using only F^2 and F^3 . The performance of this reduced ensemble is then compared to the full ensemble forecast $F^{\{1,2,3\}}$, which incorporates all three models.

On the other hand, LASOMO involves ensemble constructions from all possible subsets of models. For each subset S that does not contain the model i , $S \cup \{i\}$ plays a role of A in the LOMO; the score associated with the subset S is the difference of measures between F^S and $F^{S \cup \{i\}}$. Then, all scores are aggregated across all possible subsets that the model i does not belong to. For example, using the earlier setup of three forecast models, LASOMO considers three subsets, $\{2\}$, $\{3\}$, and $\{2, 3\}$, to calculate the importance score of model 1 (excluding all subsets that include model 1). The ensemble forecasts $F^{\{2\}}$, $F^{\{3\}}$, and $F^{\{2,3\}}$ are then compared to $F^{\{1,2\}}$, $F^{\{1,3\}}$, and $F^{\{1,2,3\}}$, respectively. The performance differences attributable to model 1's inclusion are aggregated, which results in the importance score of model 1. We note that the subsets may have different weights during the aggregating process. The **modelimportance** package considers two cases: all subsets may have uniform weights, or weights may be assigned based on the size of the subsets, similar to the concept of Shapley values. This flexibility allows for different approaches to account for the contribution of each model depending on the desired evaluation framework.

Algorithms 1 and 2 outline the steps to implement LOMO and LASOMO for a single prediction task, respectively.

Algorithm 2 Importance score using leave all subsets of models out algorithm

Input:

- Set of n individual models, $A = \{1, 2, \dots, n\}$, and their forecasts $\{F^1, F^2, \dots, F^n\}$
- Truth value, y
- Function g to build an ensemble
- Scoring rule ψ to evaluate forecast skills

Output: Importance metric of model i , $\phi^{i, \text{lasomo}}$

```

1: for  $i = 1$  to  $n$  do
2:    $\phi^{i, \text{lasomo}} \leftarrow 0$ 
3:   Make a list of non-empty subsets of  $A$  that does not contain  $i$ :  $S_1, S_2, \dots, S_{2^{n-1}-1}$ 
4:   for  $j = 1$  to  $2^{n-1} - 1$  do
5:     Assign a weight to the subset  $S_j$ :
6:     if all subsets have uniform weights then
7:        $\gamma_{S_j} \leftarrow \frac{1}{2^{n-1} - 1}$ 
8:     else (subset's weight depends on its size)
9:        $\gamma_{S_j} \leftarrow \frac{1}{(n-1) \binom{n-1}{|S_j|}}$ 
10:    end if
11:     $F^{S_j} \leftarrow g(\{F^j \mid j \in S_j\})$ 
12:     $F^{S_j \cup \{i\}} \leftarrow g(\{F^i, F^j \mid j \in S_j\})$ 
13:    Compute  $\psi(F^{S_j}, y)$  and  $\psi(F^{S_j \cup \{i\}}, y)$ 
14:     $\phi^{i, \text{lasomo}} \leftarrow \phi^{i, \text{lasomo}} + \gamma_{S_j} \times [\psi(F^{S_j}, y) - \psi(F^{S_j \cup \{i\}}, y)]$ 
15:  end for
16: end for

```

4. Evaluating ensemble members

...

5. Examples

...

6. Summary and discussion

Multi-model ensemble forecasts often provide better accuracy and robustness than single models, and are widely used in decision-making and policy planning across a variety of domains. The contribution of each component model to the accuracy of the ensemble depends on its own unique characteristics. The **modelimportance** package enables the quantification of the value that each component model adds to the ensemble performance in different evaluation

contexts.

References

- Bosse NI, Gruson H, Cori A, van Leeuwen E, Funk S, Abbott S (2022). “Evaluating forecasts with scoringutils in R.” *arXiv preprint arXiv:2205.07090*.
- Bracher J, Ray EL, Gneiting T, Reich NG (2021). “Evaluating epidemic forecasts in an interval format.” *PLOS computational biology*, **17**(2), e1008618. doi:10.1371/journal.pcbi.1008618.
- Chen T, He T, Benesty M, Khotilovich V, Tang Y, Cho H, Chen K, Mitchell R, Cano I, Zhou T, Li M, Xie J, Lin M, Geng Y, Li Y, Yuan J (2024). *xgboost: Extreme Gradient Boosting*. R package version 1.7.7.1, URL <https://CRAN.R-project.org/package=xgboost>.
- Consortium of Infectious Disease Modeling Hubs (2024). “The hubverse: open tools for collaborative forecasting.” Accessed: 2025-01-14, URL <https://hubverse.io/en/latest/index.html>.
- Gneiting T, Raftery AE (2005). “Weather forecasting with ensemble methods.” *Science*, **310**(5746), 248–249.
- Gneiting T, Raftery AE (2007). “Strictly Proper Scoring Rules, Prediction, and Estimation.” *Journal of the American Statistical Association*, **102**(477), 359–378. doi:10.1198/016214506000001437.
- Hamner B, Frasco M (2018). *Metrics: Evaluation Metrics for Machine Learning*. R package version 0.1.4, URL <https://CRAN.R-project.org/package=Metrics>.
- Hastie T, Tibshirani R, Friedman J (2001). *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA.
- He K, Yang Q, Ji L, Pan J, Zou Y (2023). “Financial Time Series Forecasting with the Deep Learning Ensemble Model.” *Mathematics*, **11**(4). doi:10.3390/math11041054.
- Hyndman RJ, Khandakar Y (2008). “Automatic time series forecasting: the forecast package for R.” *Journal of Statistical Software*, **27**(3), 1–22. doi:10.18637/jss.v027.i03.
- Jordan A, Krüger F, Lerch S (2019). “Evaluating Probabilistic Forecasts with scoringRules.” *Journal of Statistical Software*, **90**(12), 1–37. doi:10.18637/jss.v090.i12.
- Jordan A Guerra, Sophie A Murray, D Shaun Bloomfield, Peter T Gallagher (2020). “Ensemble forecasting of major solar flares: methods for combining models.” *J. Space Weather Space Clim.*, **10**, 38. doi:10.1051/swsc/2020042.
- Kim M, Ray EL, Reich NG (2024). “Beyond forecast leaderboards: Measuring individual model importance based on contribution to ensemble accuracy.” URL <https://arxiv.org/abs/2412.08916>.
- Kuhn M (2008). “Building Predictive Models in R Using the caret Package.” *Journal of Statistical Software*, **28**(5), 1–26. doi:10.18637/jss.v028.i05. URL <https://www.jstatsoft.org/index.php/jss/article/view/v028i05>.

- Liaw A, Wiener M (2002). “Classification and Regression by randomForest.” *R News*, **2**(3), 18–22. URL <https://CRAN.R-project.org/doc/Rnews/>.
- Lutz CS, Huynh MP, Schroeder M, Anyatonwu S, Dahlgren FS, Danyluk G, Fernandez D, Greene SK, Kipshidze N, Liu L, Mgbere O, McHugh LA, Myers JF, Siniscalchi A, Sullivan AD, West N, Johansson MA, Biggerstaff M (2019). “Applying infectious disease forecasting to public health: a path forward using influenza forecasting examples.” *BMC Public Health*, **19**(1), 1659. doi:10.1186/s12889-019-7966-8.
- Ray EL, Reich NG (2018). “Prediction of infectious disease epidemics via weighted density ensembles.” *PLOS computational biology*, **14**(2), e1005910. doi:10.1371/journal.pcbi.1005910.
- Reich NG, McGowan CJ, Yamana TK, Tushar A, Ray EL, Osthus D, Kandula S, Brooks LC, Crawford-Crudell W, Gibson GC, Moore E, Silva R, Biggerstaff M, Johansson MA, Rosenfeld R, Shaman J (2019). “Accuracy of real-time multi-model ensemble forecasts for seasonal influenza in the U.S.” *PLOS computational biology*, **15**(11), e1007486. doi:10.1371/journal.pcbi.1007486.
- Ridgeway G, Developers G (2024). *gbm: Generalized Boosted Regression Models*. R package version 2.2.2, URL <https://CRAN.R-project.org/package=gbm>.
- Shandross L, Howerton E, Contamin L, Hochheiser H, Krystalli A, of Infectious Disease Modeling Hubs C, Reich NG, Ray EL (2024). “hubEnsembles: Ensembling Methods in R.” *medRxiv*. doi:10.1101/2024.06.24.24309416. <https://www.medrxiv.org/content/early/2024/06/25/2024.06.24.24309416.full.pdf>.
- Sun S, Wang S, Wei Y (2020). “A new ensemble deep learning approach for exchange rates forecasting and trading.” *Advanced Engineering Informatics*, **46**, 101160. doi:<https://doi.org/10.1016/j.aei.2020.101160>.
- Viboud C, Sun K, Gaffey R, Ajelli M, Fumanelli L, Merler S, Zhang Q, Chowell G, Simonsen L, Vespignani A (2018). “The RAPIDD ebola forecasting challenge: Synthesis and lessons learnt.” *Epidemics*, **22**, 13–21. doi:10.1016/j.epidem.2017.08.002.
- Yan Y (2024). *MLmetrics: Machine Learning Evaluation Metrics*. R package version 1.1.3, URL <https://CRAN.R-project.org/package=MLmetrics>.

Affiliation:

Minsu Kim

E-mail: minsukim@umass.edu

Nicholas G. Reich