



ARCHIVOS DE Bronconeumología

www.archbronconeumol.org



Discussion Letter

Reevaluating Feature Selection in Machine Learning Models for Identifying Disease-Modifying Agents in Obstructive Sleep Apnea

To the Director,

Belmonte et al. conducted a comprehensive study on the synergistic integration of the miRNome, machine learning, and bioinformatics to identify potential disease-modifying agents for obstructive sleep apnea.¹ Their methodology involved variable selection utilizing random forests and sparse partial least squares. To enhance the robustness and consistency of their selection process, they implemented a rigorous approach by repeating each method 50 times.¹

While Belmonte et al. showcased innovative machine learning models for identifying potential disease-modifying agents for obstructive sleep apnea, this paper raises critical concerns regarding the reliance on feature selection from random forests due to the model-specific nature of feature importance. It is vital for Belmonte et al. to recognize that feature importance metrics derived from machine learning models can be inherently biased, as different models employ distinct methodologies for calculating these values due to the absence of ground truth values for accuracy validation.^{2,3}

Random forests, as a supervised machine learning approach, rely on ground truth values for validating prediction accuracy. However, the validation does not extend to the reliability of feature importance. High prediction accuracy does not necessarily indicate trustworthy feature importances; thus, caution is warranted when interpreting these metrics. Over 100 peer-reviewed articles discussed non-negligible biases in feature importances from models.^{2,3}

To accurately evaluate genuine associations between the target variable and the features, three critical elements must be considered: the distribution of the data, the statistical relationships among the variables, and the validation of statistical significance through p-values. Consequently, the choice between parametric and nonparametric methods, as well as linear and nonlinear approaches, plays a pivotal role in ensuring a robust and precise analysis.

In this context, the paper recommends the use of nonlinear and nonparametric robust statistical methods,⁴ such as Spearman's correlation and Kendall's tau, supplemented by p-values for significance testing. Additionally, before applying these statistical methods, it is essential to assess multicollinearity by executing the Variance Inflation Factor (VIF) analysis, which can help to eliminate features with collinearity and interactions, thereby reducing feature inflation.⁵

While Belmonte et al. made significant strides in integrating the miRNome, machine learning, and bioinformatics for identifying disease-modifying agents in obstructive sleep apnea, it

is crucial to address concerns regarding the reliability of feature selection from random forests. Acknowledging the inherent biases in feature importance metrics derived from machine learning models—especially in the absence of ground truth values for validation—is vital for producing credible results. To ensure accurate associations between the target variable and features, careful consideration of data distribution, variable relationships, and statistical significance is essential. By employing robust statistical methods, such as Spearman's correlation and Kendall's tau, alongside VIF analysis for multicollinearity, they can enhance the rigor and validity of their findings in this evolving field.

Authors' Contributions

Yoshiyasu Takefuji completed this research and wrote this article.

Ethics Approval

Not applicable.

Consent to Participate

Not applicable.

Consent for Publication

Not applicable.

Artificial Intelligence Involvement

Not applicable.

Funding

This research has no fund.

Conflicts of Interest

The author has no conflict of interest.

Availability of Data and Material

Not applicable.

Code Availability

Not applicable.

<https://doi.org/10.1016/j.arbres.2025.01.014>

0300-2896/© 2025 SEPAR. Published by Elsevier España, S.L.U. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

Please cite this article as: Takefuji Y, Reevaluating Feature Selection in Machine Learning Models for Identifying Disease-Modifying Agents in Obstructive Sleep Apnea, Archivos de Bronconeumología, <https://doi.org/10.1016/j.arbres.2025.01.014>

References

1. Belmonte T, Benitez ID, García-Hidalgo MC, Molinero M, Pinilla L, Mínguez O, et al. Synergic integration of the miRNome, machine learning and bioinformatics for the identification of potential disease-modifying agents in obstructive sleep apnea. *Arch Bronconeumol*. 2024, <http://dx.doi.org/10.1016/j.arbres.2024.11.011>. ISSN 0300-2896.
2. Fisher A, Rudin C, Dominici F. All models are wrong, but many are useful: learning a variable's importance by studying an entire class of prediction models simultaneously. *J Mach Learn Res*. 2019;20:177.
3. Strobl C, Boulesteix AL, Zeileis A, Hothorn T. Bias in random forest variable importance measures: illustrations, sources and a solution. *BMC Bioinformatics*. 2007;8:25, <http://dx.doi.org/10.1186/1471-2105-8-25>.
4. Okoye K, Hosseini S. Correlation tests in R: Pearson Cor, Kendall's Tau, and Spearman's Rho R programming. Singapore: Springer; 2024, http://dx.doi.org/10.1007/978-981-97-3385-9_12.
5. Salmerón-Gómez R, García-García CB, García-Pérez J. A redefined variance inflation factor: overcoming the limitations of the variance inflation factor. *Comput Econ*. 2025;65:337–63, <http://dx.doi.org/10.1007/s10614-024-10575-8>.

Yoshiyasu Takefuji

Faculty of Data Science, Musashino University, Tokyo, Japan

E-mail address: takefuji@keio.jp