
COMMON PITFALLS IN EVALUATING MODEL PERFORMANCE AND STRATEGIES FOR AVOIDANCE IN AGRICULTURAL STUDIES

A PREPRINT

 **C. P. James Chen***

School of Animal Sciences
Virginia Tech
Blacksburg, VA 24061
niche@vt.edu

 **Robin R. White**

School of Animal Sciences
Virginia Tech
Blacksburg, VA 24061
rrwhite@vt.edu

Ryan Wright

School of Animal Sciences
Virginia Tech
Blacksburg, VA 24061
ryanw22@vt.edu

January 4, 2025

ABSTRACT

Predictive modeling is a cornerstone of data-driven research and decision-making in precision agriculture, yet achieving robust, interpretable, and reproducible model evaluations remains challenging. This study addresses two central issues in model evaluation — methodological pitfalls in cross-validation (CV) and data-structure effects on performance metrics — across five simulation experiments supplemented by real-world data. First, we show how the choice of estimator (e.g., K-fold, LOOCV) and sample size affects the reliability of performance estimates: although LOOCV can be unbiased for error-based metrics, it systematically underestimates correlation-based metrics. Second, we demonstrate that reusing the test data during model selection (e.g., feature selection, hyperparameter tuning) inflates performance estimates, reinforcing the need for proper separation of training, validation, and test sets. Third, we reveal how ignoring experimental block effects, such as seasonal or herd variations, introduces an upward bias in performance measures highlighting the importance of block CV when predictions are intended for new, unseen environment. Fourth, we highlight that different regression metrics — ranging from correlation-based (e.g., r , Lin's CCC) to error-based (e.g., RMSE, MAE) — capture distinct aspects of predictive performance and under varying error biases and variances. Finally, for classification tasks, class imbalance and threshold settings significantly alter performance metrics, illustrating why a single metric rarely suffices to characterize model performance comprehensively. Collectively, these findings emphasize the need for

*Corresponding author: James Chen <niche@vt.edu>

18 careful alignment between modeling objectives, CV strategies, and metric selection, thereby ensuring
19 trustworthy and generalizable model assessments in precision agriculture and beyond.

20 **Keywords** Model Evaluation · Performance Metrics · Simulation Studies

21 **1 Introduction**

22 **1.1 Modeling**

23 Modeling is an essential tool for hypothesis formulation and decision-making. It functions as a structured investigatory
24 framework that allows researchers to explore system understanding through the summary and analysis of empirical data.
25 Carefully constructed and evaluated models offer the potential to extend this understanding by enabling the extrapolation
26 of results to novel trials and conditions. Although only one focus of the science of modeling, the predictive role is
27 often explicitly or implicitly the ultimate goal of models derived within the precision agriculture context. Through
28 this lens, modeling provides opportunity to standardize and formalize research advancement, through developing
29 quantitative constructs that accumulate prior knowledge derived by the broader the scientific community. Evaluating
30 model performance becomes particularly critical when considering this role within the knowledge generation enterprise,
31 necessitating a rigorous and standardized approach that allows for both reproducibility and comparability. As more and
32 more model-based exercises are developed using slightly different methods, or slightly different datasets, it becomes
33 increasingly challenging to evaluate, characterize, compare, and balance information generated by the resulting modeling
34 tools, particularly when results are conflicting. Specifically, reporting model performance through poorly-defined
35 metrics or incomplete procedures can create opportunity for confusion, misinterpretation, and miscommunication, and
36 can ultimately result in distrust in model-based tools and impede scientific progress.

37 This study examines two primary challenges that arise during model evaluation: those associated with the evaluation
38 methodology and those stemming from the data structure. The former emphasizes the reliability of estimated perfor-
39 mance and essential measures to avoid overestimating a model's capabilities. The latter depends on the nature of the
40 modeling exercise: for regression tasks, variance and bias are particularly important for assessing performance, whereas
41 for classification tasks, class imbalance poses a critical concern. Employing multiple performance metrics can help
42 prevent misinterpretation due to these factors. To illustrate the significance of these challenges and effective strategies
43 to address them, we conduct a series of simulations complemented by real-world data examples.

44 **1.2 Model Evaluation**

45 Model evaluation in the context of predictive analytics seeks to explore how well a model can generalize to new
46 prediction contexts not seen during model training. Although commonly referred to as "model validation" in the
47 literature, this term implies a false degree of confidence given that the word "validation" means to prove something
48 true. There is no single test, or recognized suite of tests, to prove a model valid. Instead, the term "evaluation," which
49 involves assessing the value, nature, character, or quality of something, is more fitting. It is essential to evaluate model
50 performance on unseen data to ensure the approach is applicable to new experiments. To this end, cross-validation (CV)
51 is widely recognized as a standard method for model evaluation.

52 The most common CV method is K-fold CV, which partitions the dataset into K equally sized folds. In each iteration,
53 one fold is reserved as the test set (i.e., new data, noted as $\mathcal{D}_{\text{test}}$), while the remaining folds are used as the training set

54 (noted as $\mathcal{D}_{\text{train}}$) to construct the model. Once the model is trained, it is evaluated on the $\mathcal{D}_{\text{test}}$ to obtain an estimate of
55 the model performance \hat{g} . The process will iterate K times until each fold has been used as the $\mathcal{D}_{\text{test}}$ once. The average
56 performance over all K folds is deemed as the expected generalization performance of the model $\mathbb{E}[\hat{g}]$ on new data.

57 However, there is always an evaluation bias between the estimated performance $\mathbb{E}[\hat{g}]$ and the true generalization
58 performance G , which can only be approximated by evaluating the same model on an infinite number of unseen data.
59 Depending on the performance metric used in evaluation, a positive evaluation bias ($\mathbb{E}[\hat{g}] - G$) typically suggests that the
60 model evaluation procedure concludes a pessimistic estimation of the model performance, since the true performance
61 is expected to be lower than the estimated performance. Another aspect of model evaluation error is the variance of
62 each estimated performance \hat{g} across the K folds. For example, there are five estimates in a 5-fold cross-validation.
63 The variance among these five estimates is defined as the evaluation variance. A high evaluation variance suggests that
64 the performance is sensitive to the choice of data folds, and a small size or an over-complex model can lead to a high
65 evaluation variance.

66 There is a trade-off relationship between evaluation bias and variance, which can be understood through the framework
67 of the squared evaluation bias (see Appendix for a detailed derivation). When performing K-fold CV with a fixed
68 sample size and model complexity, the choice of K is the pivotal element shaping the model evaluation. When the K is
69 set to a larger value; each training set $\mathcal{D}_{\text{train}}$ is larger in size, resulting in a model trained on a more representative subset
70 of the population of interest, leading to lower bias. However, a large K comes with a trade-off: the corresponding test
71 subset $\mathcal{D}_{\text{test}}$ is compressed in size, making the tested model more sensitive to the specific data points, and thus inflating
72 the validation variance. Conversely, a smaller K, along with a minor training set $\mathcal{D}_{\text{train}}$, reduces their representativeness
73 and increases bias. Nevertheless, a larger size of the test set $\mathcal{D}_{\text{test}}$ leads to more consistent estimations across the folds
74 and, consequently, reduces the validation variance.

75 Leave-one-out cross-validation (LOOCV) is a variant of K-fold CV where K equals the sample size of the complete
76 dataset \mathcal{D} . It provides an unbiased estimation of model performance because the training set $\mathcal{D}_{\text{train}}$ closely resembles the
77 unseen population of interest, given its size of $N - 1$, where N is the sample size. However, as the trade-off discussion
78 suggested, this method can lead to high validation variance due to the model being evaluated on one sample at a time.
79 The true unbiased nature of LOOCV is fully realized only when all K folds are utilized. Performing an incomplete
80 LOOCV can introduce significant bias because of the inherent high validation variance, which often occurs when
81 training each model iteration is prohibitively time-consuming or computationally demanding. In specific contexts, such
82 as genomic prediction, strategies like the one described by Cheng et al. leverage the matrix inverse lemma, which
83 allows for computational savings by avoiding the inversion of large matrices in each fold. This technique significantly
84 reduces the dependency of computational resources on the sample size [1]. Van Dixhoorn et al. exemplify the use of
85 LOOCV with a small dataset, aiming to predict cow resilience with limited data resources [2]. Nevertheless, for large
86 datasets, LOOCV is generally not recommended due to computational inefficiency. Further details of bias-variance
87 trade-off have been extensively explored in the statistical literature [3, 4].

88 1.3 Model Selection

89 Model selection becomes necessary when models are not entirely determined by the data alone. For example, in a
90 regularized linear regression model such as a ridge regression [5] or the least absolute shrinkage and selection operator
91 (LASSO) [6], it is essential to define a regularization parameter, λ , before fitting the model to the data. A larger λ value
92 yields a more regularized model, which tends to reduce smaller coefficients to negligible values or zero. This approach
93 helps in preventing overfitting noise in the training data. The definition of loss functions for the regularized models
94 were described in S.9 and S.10 of the Appendix.

95 These pre-defined parameters, like λ , influence model fitting and remain constant during the training process. Such
96 parameters are referred to as hyperparameters. Beyond regularized models, hyperparameters are crucial in other
97 predictive models, enhancing flexibility and robustness. For example, in the Support Vector Regression (SVR) [7],
98 the regressors X are projected onto a linear subspace to approximate the target variable Y . By choosing a suitable
99 kernel function, which transforms the regressors into a non-linear space, as a hyperparameter, SVR can more effectively
100 capture non-linear relationships, thus significantly improving model performance. Another hyperparameter example is
101 the number of latent variables in the Partial Least Square (PLS) Regression [8], which condenses the original regressors
102 into a more manageable set of latent variables, reducing multicollinearity issues. Fewer latent variables might lose
103 significant information from the original regressors, while too many can lead to overfitting. Similarly, in Random
104 Forest [9], hyperparameters such as tree depth and the number of trees influence model complexity by dictating how
105 many feature splits are possible and how many weak learners comprise the ensemble. The same principle applies to
106 convolutional neural networks, where increasing the number of hidden layers or filter sizes can capture more complex
107 patterns in the data but also heightens the risk of overfitting [10]. All these examples highlight the fact that selecting the
108 most suitable hyperparameters, which is known as hyperparameter tuning, is crucial for optimizing model performance.
109 Feature selection is another crucial aspect of model selection. This process involves fitting the model to a selected
110 subset of the original features, particularly essential in high-dimensional data scenarios where the number of features
111 exceeds the number of observations, leading to poor model generalization. For instance, Ghaffari et al. sought to predict
112 health traits in 38 multiparous Holstein cows using a metabolite profiling strategy. Out of 170 metabolites, only 12
113 were identified as effective discriminators between healthy and over-conditioned cows and were thus selected for the
114 predictive model [11]. Therefore, optimizing feature subsets is a vital model selection strategy that significantly affects
115 model performance. Including the model selection process within the cross-validation is essential to avoid common
116 pitfalls. The risk of inflated model performance arises when model selection is guided by results on the test dataset.
117 Even if the chosen model is subjected to k-fold cross-validation afterward, its selection bias toward the test set can
118 lead to overestimating its efficacy. This issue has been highlighted in statistical literature [3]. A practical solution is to
119 divide the dataset into training, validation, and test sets. The validation set is then used for model selection, ensuring the
120 test set remains completely unused during the training phase, thereby providing a more accurate measure of model
121 performance. For instance, the study by Rovere et al. exemplifies best practices in hyperparameter tuning and feature
122 selection by employing an independent cross-validation step prior to assessing model performance. This approach

123 enabled the precise selection of relevant spectral bands from the mid-infrared spectrum and the optimal number of
124 latent dimensions in PLS with Bayesian regression for predicting the fatty acid profile in milk [12]. Similarly, Becker et
125 al. demonstrated a robust evaluation by using nested cross-validation loops; the inner loop conducted a grid search
126 for the best hyperparameters in logistic regression, while the outer loop was designed to evaluate the performance
127 of the resulting optimized model [13]. Both examples underscore the importance of separating model selection from
128 performance evaluation to ensure the validity and reliability of the results.

129 **1.4 Cross Validation Design with Block Effects**

130 Blocking is an essential approach in experimental design to control for variations that can confound the variable of
131 interest. For instance, Lahart et al. investigated the dry matter intake of grazing cows using mid-infrared (MIR)
132 spectroscopy technology across multiple herds under varying experimental conditions [14]. Given the significant
133 variation between herds, which may contribute to individual differences in both dry matter intake (i.e., response variable)
134 and MIR spectra (i.e., independent variables), it is crucial to consider the herd as a blocking factor before evaluating the
135 predictability of dry matter intake using MIR spectra. This consideration should also extend to model evaluation. In the
136 cited study, variations in dry matter intake, the primary focus of the prediction model, were observed to exceed one
137 standard deviation among some herds. In cross-validation, if samples from the same herd are assigned to different folds,
138 with one fold used as the test set, the model is likely to achieve high accuracy. This accuracy may largely result from
139 explaining the inter-herd variation rather than individual variations in dry matter intake, leading to an overestimation of
140 model performance. To avoid this pitfall, block cross-validation, where each block (i.e., herd in this example) is used as
141 a fold, is recommended for unbiased model evaluation. Literature reviews have indicated that block cross-validation
142 effectively evaluates model performance on external or unseen datasets [15]. In the same study by Lahart et al., three
143 cross-validation strategies were compared: random cross-validation (Random CV), which randomly assigns samples
144 to folds; within-herd validation, training and testing the model within each herd; and across-herd validation (Block
145 CV), where each herd is used as a fold and tested in turn. The results showed that performance estimates in block CV
146 were noticeably lower than the other two strategies, supporting the hypothesis that ignoring block effects inflates model
147 performance. Other studies considering block effects, including diet [16], herd [12], and farm location [17, 18], have
148 shown similar results in cross-validation, demonstrating block CV's effectiveness in evaluating model performance on
149 external datasets.

150 **1.5 Model Performance Metrics**

151 Model performance metrics serve as quantitative indicators for evaluating model performance. They are critical for
152 benchmarking various modeling approaches and for evaluating hypotheses underpinning these different approaches.
153 Choosing appropriate metrics to support hypothesis testing is crucial, as in-ideal selection may lead to overly optimistic
154 conclusions. Due to the different goals of regression and classification tasks, it is critical to ensure that these different

155 model types are evaluated using different metrics. As such, metrics for regression and classification are discussed
 156 individually.

157 **1.5.1 Metrics in Regression Tasks**

Table 1: Summary of model performance metrics for regression tasks.

| Metric | Type | Scale-invariant | Range |
|---|-----------------|-----------------|----------------|
| Root mean square error (RMSE) | Error-based | No | $[0, \infty]$ |
| Mean absolute error (MAE) | Error-based | No | $[0, \infty]$ |
| Root mean squared percentage error (RMSPE) | Error-based | Yes | $[0, \infty]$ |
| Root mean standard deviation ratio (RSR) | Error-based | Yes | $[0, \infty]$ |
| Pearson's correlation coefficient (r) | Linearity-based | Yes | $[-1, 1]$ |
| Coefficient of determination (R^2) | Linearity-based | Yes | $[-\infty, 1]$ |
| Lin's concordance correlation coefficient (CCC) | Linearity-based | Yes | $[-1, 1]$ |

158 Regression models aim to predict continuous variables and are commonly employed in diverse applications, such as
 159 estimating body condition scores [19, 20], body weight [21, 22], milk composition [12, 18, 23, 24], efficiency of feed
 160 resource usage [16, 25, 26], and early-lactation behavior [2]. The metrics in regression tasks evaluate the agreement
 161 between the predicted value \hat{y} and the true values y . The agreement can be generally quantified in two ways: error-based
 162 metrics and linearity-based metrics. The metrics are summarized in Table 1.

163 Error-based metrics focus on the deviation of each pair of predicted and true values, while linearity-based metrics
 164 consider overall linear relationships between the predictions and the truths. The root mean square error (RMSE) and the
 165 mean absolute error (MAE) are two common error-based metrics:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (1.1)$$

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (1.2)$$

166 where y_i and \hat{y}_i are the true and predicted values, respectively, and n is the sample size. Both metrics preserve the scale
 167 of the original data, making them easy to interpret in real-world units. Additionally, compared to MAE, RMSE penalizes
 168 large errors more due to the squared term, making it more sensitive to outliers. In the cow production, monitoring
 169 animal body weight is a common practice to aid in the management of dairy cows. Studies by Song et al. and Xavier et
 170 al. have utilized RMSE to assess the effectiveness of three-dimensional cameras in estimating dairy cow body weight,
 171 yielding RMSE values of 41.2 kg and 12.1 kg, respectively [21, 22]. These figures provide a straightforward value for
 172 farmers to gauge whether the prediction error is tolerable, considering their specific operational costs and management
 173 thresholds. In essence, RMSE translates complex model accuracy into practical insights for productive agricultural
 174 units. When evaluating the same model across different traits, which may have different scales, a common practice is to
 175 express error metrics in a scale-free manner. This can be achieved by expressing RMSE as a percent of the deviation

176 from the observed value, such as root mean squared percentage error (RMSPE), or as a Root Mean Standard Deviation
 177 Ratio (RSR) that normalizes the RMSE by the standard deviation of the observed values:

$$\text{RMSPE} = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{y_i} \right)^2} \quad (1.3)$$

$$\text{RSR} = \frac{\text{RMSE}}{\sigma_y} \quad (1.4)$$

178 where σ_y is the standard deviation of the observed values. When expressed as a percent, RMSPE typically ranges from
 179 0 and above, with values closer to 0 indicating perfect prediction. Much like expressing RMSE as a percent, RSR is
 180 valuable to interpret RMSE in terms of the context of the variance in the observations. Values below 1 suggest that the
 181 model yields predictions less variable than the standard deviation, while values above 1 suggest that the prediction is
 182 imprecise.

183 On the other hand, Pearson's correlation coefficients (r) and the coefficient of determination (R^2) are two common
 184 linearity-based metrics:

$$\begin{aligned} r &= \frac{\text{cov}(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}} \\ &= \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2 \sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}} \end{aligned} \quad (1.5)$$

$$\begin{aligned} R^2 &= 1 - \frac{SS_{\text{residual}}}{SS_{\text{total}}} \\ &= 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \end{aligned} \quad (1.6)$$

185 where SS_{residual} is the residual sum of squares and SS_{total} is the total sum of squares. Each y_i and \hat{y}_i are the i th elements
 186 of the actual response vector y and the predicted response vector \hat{y} , respectively. \bar{y} and $\bar{\hat{y}}$ are their respective means.
 187 Both r^2 and R^2 are scale invariant, meaning their values are unaffected by the scale of the observed data because they
 188 are normalized by the variation in the denominator.

189 The correlation coefficient r measures the strength of the linear relationship between two continuous variables, y and \hat{y} ,
 190 and ranges from -1 to 1. A value of 0 indicates no prediction accuracy in the evaluated model. One special characteristic
 191 of correlation r is that it is unaffected by the scale of the predictions or biases; it focuses on the relative changes
 192 in the predicted values compared to the true values. Thus, even if the prediction biases are scaled up or down, the
 193 correlation r between \hat{y} and y remains the same. This property is particularly useful when the focus is more on ranking
 194 predictions rather than their absolute values. For example, this metric has been used to evaluate models that identify
 195 high-performing production individuals, demonstrating the ability to predict nutrient digestibility in dairy cows [26] and
 196 to select models based on their ability to rank traits such as feed intake and milk composition in dairy cows [27, 12].

197 The coefficient of determination R^2 quantifies model performance from the proportion of variance in the dependent
 198 variable that is predictable from the independent variables. It ranges from negative infinity to 1, where 1 indicates
 199 that the model explains all the variance in the dependent variable, and 0 indicates that the model performs no better
 200 than predicting all samples as the mean of the observed values. R^2 is useful in comparing multiple regression models,
 201 as demonstrated in studies that regress body weight of dairy cows on a set of morphological traits [22], examine
 202 the relationship between milk spectral profiles and nitrogen utilization efficiency [16], and evaluate the predictive
 203 performance of milk fatty acid composition [23].

204 It worth noting that many literatures have misinterpreted the relationship between r and R^2 . The coefficient of
 205 determination R^2 is not always equivalent to the square of the correlation coefficient r^2 . The equivalence only holds
 206 when the same dataset is used for both model fitting and evaluation in a least squares regression model. The model
 207 assumes a zero covariance between the fitted residual and the predicted values \hat{y} , and it also assumes that the residuals
 208 (i.e., prediction biases) are centered on zero. In practice when predictions are made on new data, those assumptions
 209 are often violated, leading to discrepancies between r^2 and R^2 . A details derivation of the equivalence is provided in
 210 Equation S.11 S.12 in the Appendix.

211 In addition to r^2 and R^2 , another linearity-based metric is Lin's concordance correlation coefficient (CCC) [28]:

$$\begin{aligned} \text{CCC} &= \frac{2r\sigma_y\sigma_{\hat{y}}}{\sigma_y^2 + \sigma_{\hat{y}}^2 + (\bar{y} - \hat{\bar{y}})^2} \\ &= \frac{2\text{cov}(y, \hat{y})}{\sigma_y^2 + \sigma_{\hat{y}}^2 + (\bar{y} - \hat{\bar{y}})^2} \end{aligned} \quad (1.7)$$

212 where r is the Pearson correlation coefficient. The CCC is a comprehensive metric because it considers both the
 213 correlation and the scale bias between the predicted and true values. It fills the gap left by r^2 where the scale bias is
 214 ignored. Geometrically, CCC measures how well the predicted values \hat{y} fall on the 45-degree line in a scatter plot of
 215 the predicted (x-axis) and true values (y-axis). It is advantageous over R^2 because it consistently ranges from -1 to 1,
 216 making it easier to interpret and compare across different studies. The CCC is crucial when precise predictions are
 217 required for both the scale and the rank of the trait of interest, such as in studies predicting cotton crop yields based on
 218 soil and terrain profiles [29].

219 1.5.2 Metrics in Classification Tasks

220 Classification models aim to predict categorical outcomes such as 'healthy' or 'sick,' 'susceptible' or 'resistant,' and
 221 'high yield' or 'low yield.' To evaluate classification performance, one must first establish a confidence threshold to
 222 dichotomize the prediction probabilities. For instance, if a classification model predict a sample as 'sick' with a 0.7
 223 probability, and the threshold is set at 0.5. Since the 0.7 prediction probability exceeds the threshold, the sample is
 224 predicted as a positive sample. It is worth mentioning that this threshold is adjustable to fine-tune model performance
 225 for particular focus, such as minimizing false positives or false negatives. All classification metrics are derived from

226 the confusion matrix, which summarizes the model’s performance in a 2x2 table, where the rows represent the actual
 227 classes and the columns represent the predicted classes.

Table 2: Confusion matrix for binary classification.

| | | Predicted | |
|--------|----------|---------------------|---------------------|
| | | Positive | Negative |
| Actual | Positive | True Positive (TP) | False Negative (FN) |
| | Negative | False Positive (FP) | True Negative (TN) |

228 The confusion matrix (Table 2) consists of four components: true positives (TP), true negatives (TN), false positives
 229 (FP), and false negatives (FN). Most common metrics used in classification tasks are summarized in Table 3.

Table 3: Summary of model performance metrics for classification tasks.

| Metric | Denominator | Focus |
|---------------------------|---------------------|-------------|
| True positive rate (TPR) | Actual positives | Correctness |
| True negative rate (TNR) | Actual negatives | Correctness |
| False negative rate (FNR) | Actual positives | Error |
| False positive rate (FPR) | Actual negatives | Error |
| Sensitivity | Actual positives | Correctness |
| Specificity | Actual negatives | Correctness |
| Precision | Predicted positives | Correctness |
| Recall | Actual positives | Correctness |
| Accuracy | All samples | Balance |
| F1 score | All samples | Balance |
| F-beta score | All samples | Balance |
| MCC | All samples | Balance |

230 The metrics can be characterized by two key factors: their denominator and their focus on either correctness or error.
 231 Understanding the denominator of a metric helps clarify its scope of interest. For instance, if one wants to evaluate
 232 how well the model correctly predicts positive samples, metrics that use actual positives as the denominator should be
 233 prioritized. It is noted that in Table 3, the metrics are organized in four subsections. The metrics in the first subsection
 234 have self-explanatory names, each emphasizing a specific aspect of the model’s performance:

$$\begin{aligned}
 \text{True positive rate (TPR)} &= \text{Sensitivity} \\
 &= \text{Recall} \\
 &= \frac{\text{TP}}{\text{Total Actual Positives}}
 \end{aligned} \tag{1.8}$$

$$\begin{aligned}
 \text{True negative rate (TNR)} &= \text{Specificity} \\
 &= \frac{\text{TN}}{\text{Total Actual Negatives}}
 \end{aligned} \tag{1.9}$$

235 Both TPR and TNR focus on the correctness of the model's predictions, but TPR is concerned with positive samples,
 236 while TNR is concerned with negative samples. High TPR is essential where missing a positive case has serious
 237 consequences, or where false positives are easily rectifiable. For instance, detecting lameness or abnormal gait is crucial,
 238 as these can indicate underlying pathologies [30] and impact welfare-related transport decisions [31]. An automated
 239 detection system [30, 32, 33] with high TPR can mitigate economic losses by flagging at-risk cows. The benefit here
 240 lies in the feasibility of re-examining false positives, thus preventing more severe outcomes from undetected cases.
 241 In contrast, the false negative rate (FNR) and false positive rate (FPR) focus on the model's errors:

$$\text{False negative rate (FNR)} = \frac{\text{FN}}{\text{Total Actual Positives}}$$

$$\text{False positive rate (FPR)} = \frac{\text{FP}}{\text{Total Actual Negatives}}$$

242 The second section of Table 3 includes sensitivity and specificity, which are equivalent to TPR and TNR, respectively.
 243 These terms are widely used in medical diagnostics due to their emphasis on accurately identifying true positive and
 244 true negative cases, which are critical requirement for tests and screenings for disease detection.
 245 The third section includes precision and recall, which focus on different aspects of positive cases. Machine learning
 246 community used to report precision and recall together, as the community focus more on the positive samples than the
 247 negative samples. For example, in computer vision applications, how well a model can correctly classify and localize
 248 the object of interest (positives) from an image is more important than how well the model can correctly know what area
 249 is irreleavnt background (negatives). Precision evaluates the correctness of the predicted positive cases, ensuring that
 250 the predictions are accurate, while recall measures the completeness of identifying all actual positive cases, emphasizing
 251 the model's ability to capture true positives. Precision measure the trustworthiness of positive predictions made by the
 252 model (Eq. 1.10). High precision is crucial in scenarios where false positives incur significant costs. For instance, in
 253 contexts where clinical treatments and culling are expensive, such as detecting bovine tuberculosis [34] or mastitis [35]
 254 using non-invasive methods, a high-precision model is crucial to minimize unnecessary costs and interventions from
 255 false positives. Precision and recall are a pair of metrics commonly used in machine learning applications, particularly in
 256 multi-class classification or detection scenarios. In these contexts, the evaluation of negative samples (i.e., non-positive
 257 samples) is often replaced by examining the precision and recall for each individual class. This approach allows for a
 258 more granular assessment of the model's performance across all classes, ensuring that both the quality of predictions
 259 and the ability to identify all relevant samples are accounted for.

$$\text{Precision} = \frac{\text{TP}}{\text{Total Predicted Positives}} \quad (1.10)$$

260 The last section of Table 3 includes accuracy, F1 score, F-beta score, and Matthews Correlation Coefficient (MCC).
 261 These metrics offer a balanced evaluation of the model’s performance by taking into account both correctness and error
 262 rates, as well as both positive and negative samples. Among them, accuracy is the most straightforward metric for
 263 evaluating classification models, as it measures the proportion of correctly classified samples out of the total samples.

$$\begin{aligned} \text{Accuracy} &= \frac{\text{Total Correct Predictions}}{\text{Total Predictions}} \\ &= \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \end{aligned} \quad (1.11)$$

264 It summarizes an overall model performance by calculating the proportion of correctly classified samples among all
 265 samples. Nonetheless, accuracy can be misleading when the classes are imbalanced. For example, if a study predicting
 266 the presence of a specific event, of which the prevalence was only 10%. In this case, a model that predicts all samples
 267 as negative would achieve an accuracy of 90%, which is misleadingly high. The F1 score, which is the harmonic mean
 268 of precision and recall (i.e., TPR), provides a balanced measure of model performance:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (1.12)$$

269 Unlike accuracy, the F1 score considers both false positives and false negatives by balancing precision and recall,
 270 making it a more reliable metric for imbalanced datasets. A variant of the F1 score is the F-beta score, which allows for
 271 the adjustment of the balance between precision and recall by introducing a weight parameter β :

$$\text{F-beta} = (1 + \beta^2) \times \frac{\text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}} \quad (1.13)$$

272 A common variant is the F2 score, which places more emphasis on false negatives (i.e., recall) than false positives, by
 273 setting $\beta = 2$:

$$F2 = 5 \times \frac{\text{Precision} \times \text{Recall}}{4 \times \text{Precision} + \text{Recall}} \quad (1.14)$$

274 Lastly, the Matthews correlation coefficient (MCC) considers both positive and negative samples in the dataset, providing
 275 a balanced measure of a model’s performance [36]. It is defined as:

$$MCC = \frac{\text{TP} \times \text{TN} - \text{FP} \times \text{FN}}{\sqrt{(\text{TP} + \text{FP})(\text{TP} + \text{FN})(\text{TN} + \text{FP})(\text{TN} + \text{FN})}} \quad (1.15)$$

276 The equation 1.15 symmetrically incorporates all four components of TP, TN, FP, and FN). This symmetry makes
 277 MCC invariant to class distribution changes. The coefficient ranges from -1 to 1, where 1 indicates perfect classification,
 278 0 indicates no better performance than random guessing, and -1 signifies total disagreement between prediction and

279 observation. In a case study that used feeding and daily activity behaviors to diagnose Bovine Respiratory Disease
280 in dairy calves, MCC proved particularly insightful [37]. The models in this study exhibited strong performance on
281 negative samples (i.e., healthy calves), which were more prevalent, resulting in high specificity. However, sensitivity
282 was relatively low at 0.54. In this context, MCC, with a value of 0.36, provided a more nuanced and representative
283 measure of model performance, especially given the skew towards negative samples

284 **1.6 Study Objectives**

285 This simulation study aims to highlight how biased or over-optimistic estimations of model performance usually come
286 from inappropriately conducting CV, a technique crucial for characterizing expected model performance on “new”
287 data. We demonstrate how common pitfalls, including using the exact data for both training and model assessment,
288 excluding the model selection process from CV, and neglecting experimental block effects, contribute to challenges
289 in model evaluation. Further, we scrutinize common metrics used in evaluating prediction models, including those
290 used for regression and classification tasks. Because no single metric provides a comprehensive perspective of model
291 performance, we seek, through this work, to highlight the importance of understanding the underlying theory of each
292 metric to avoid misleading conclusions.

293 There are five simulation studies being conducted to address these challenges. The first simulation study will focus
294 on the bias-variance trade-off in CV, demonstrating how the choice of K in K-fold CV affects the evaluation bias and
295 variance. The second simulation study will investigate the impact of mistakenly using the same data for model selection
296 and evaluation, highlighting the inflated model performance. The third simulation study will explore the effect of
297 excluding block effects in CV, demonstrating how ignoring block effects can lead to over-optimistic model performance.
298 The fourth simulation study will explore how various metrics respond to different combinations of bias and variance
299 in prediction errors, illustrating how these variations can lead to distinct interpretations of model performance. The
300 fifth simulation study will examine the impact of imbalanced data on classification model evaluation, highlighting
301 how the choice of metrics can influence conclusions and potentially lead to misleading interpretations. Together, this
302 series of simulation studies aims to provide guidance for researchers on accurately and consistently reporting model
303 performance, thereby promoting integrity and scientific rigor in prediction modeling research.

304 **2 Materials and Methods**

305 **2.1 Study datasets**

306 This study utilized three datasets to demonstrate the common challenges in model evaluation: A null dataset, a simulated
 307 spectral dataset, and a real-world spectral dataset.

308 **2.1.1 Null dataset**

309 The null dataset serves as a baseline for the null hypothesis, designed to evaluate the risk of introducing bias in the
 310 estimation of model performance. In this dataset, the predictors X and the target variable y are independently drawn
 311 from the same normal distribution, ensuring no linear or nonlinear relationship between the input features and the target
 312 variable:

$$\begin{cases} X \sim \mathcal{N}(0, 1) \\ y \sim \mathcal{N}(0, 1) \end{cases} \quad (2.1)$$

313 If any model evaluation exercise applied to this dataset produces a significant performance metric, it would indicate a
 314 potential bias in the evaluation process. This serves as a critical check to ensure that the evaluation methodology does
 315 not artificially inflate the perceived performance of the model.

316 **2.1.2 Simulated spectral dataset**

317 To further investigate how these identified challenges impact data with complex structures, both simulated and real
 318 spectral datasets were utilized. Spectral data is commonly encountered in agricultural studies, where the target variable
 319 is predicted using a series of spectral measurements. This type of data serves as an excellent example for this study
 320 because it often presents a significant challenge due to the strong collinearity among predictors. Effectively selecting
 321 predictors with reduced correlation is essential to mitigate overfitting and improve model robustness. The simulated
 322 spectral dataset was generated following the procedure outlined in [38], which characterizes spectral signals X as the
 323 outcome of a linear combination of a score matrix T and a loading matrix P :

$$X_{n \times m} = T_{n \times k} P_{m \times k}^\top + E_{n \times m} \quad (2.2)$$

324 where X is the spectral data matrix with n samples and m spectral variables, T is the score matrix with n samples
 325 and k latent variables, P is the loading matrix with m variables and k latent variables, and E is the residual matrix to
 326 simulate noise (Figure 1). The k latent variables represent the underlying structure of the spectral data, such as the
 327 peaks and valleys in the spectrum. In this study, the spectral data consists of 300 spectral bands ($m = 300$) and four
 328 latent variables ($k = 4$). Among these latent variables, only the first two are assumed to contribute meaningfully to the
 329 target variable, while the other two are treated as noise.

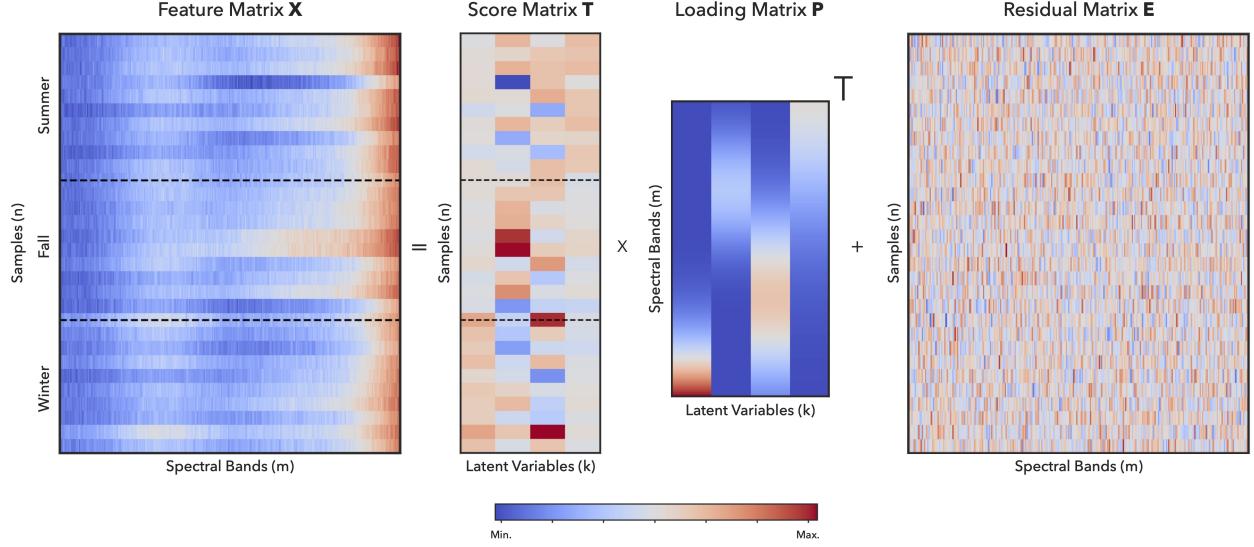


Figure 1: Matrix decomposition of the simulated spectral data. The spectral data matrix X is generated as a linear combination of the score matrix T and the loading matrix P , with added noise E . The color scale is independently normalized for each matrix.

- 330 The score matrix T defines how each sample contributes to each latent variable. For example, if the third sample
 331 exhibits higher spectral measurements around the first peak (as defined by the first latent variable), the value in the
 332 third row and first column of the score matrix will be higher relative to other rows. In this study, T was sampled from
 333 a multivariate normal distribution with a mean vector of $[1, 1, 1, 1]$ and standard deviations of $[0.02, 0.10, 0.10, 0.02]$.
 334 This setup reflects a scenario where the second and third latent variables (corresponding to specific peaks) are more
 335 pronounced compared to the first and fourth latent variables. It is inspired by the spectrum measured in the past work
 336 [39], which used 150 hyperspectral bands ranging from 1,000 nm to 1600 nm to evaluate the wheat kernel quality trait.
 337 The loading matrix P defines how each spectral variable contributes to each latent variable. Each latent variable in P
 338 was simulated using a Gaussian probability function with peaks at the -30th, 90th, 200th, and 345th spectral positions
 339 and standard deviations of $[100, 40, 60, 60]$ to simulate the width of the peaks. Negative peak positions simulate signals
 340 outside the measured spectral range. The residual matrix E is sampled from a normal distribution $\mathcal{N}(0, 0.01)$ to
 341 simulate the noise in the spectral data.
 342 Seasonal variation is an important factor in agricultural studies and is often overlooked in model evaluation. To
 343 incorporate this effect, the spectral measurements were simulated across three seasons, with random effects applied
 344 to the latent variables. These seasonal effects were modeled by multiplying different scalars with the latent variables
 345 in the score matrix T . The scalars for the first two latent variables were $[1.00, 1.10, 1.07]$, and for the latter two were
 346 $[1.07, 1.00, 1.00]$. This setup reflects a scenario where the first two latent variables are more pronounced in the second
 347 and third seasons, while the latter two latent variables dominate in the first season.
 348 The response variable y was generated as a nonlinear function of selected spectral variables. Specifically, four spectral
 349 bands (B) were selected at indices $[50, 100, 180, 230]$ from the 300 spectral bands. Nonlinear effects were introduced by

350 applying a sinusoidal transformation to the selected spectral variables, raised to the power of three (cubic nonlinearity).
 351 To ensure that each effect has a unique sinusoidal component, a phase shift was added to each effect. The phase shift is
 352 defined as $\frac{i\pi}{m}$, where i is the index of the selected spectral band, and m is the total number of selected spectral bands.
 353 This approach introduces variation in the sinusoidal behavior of each effect, ensuring they are distinct:

$$y = \sum_{i=1}^m \sin(b_i^3 + \frac{i\pi}{m}), \quad b_i \in B$$

354 where b_i represents the i -th selected spectral band from the four bands B . Finally, Gaussian noise was added to the
 355 response variable y to simulate measurement or modeling errors. The noise was generated with a standard deviation
 356 equal to that of the response variable y , simulating a scenario where only 50% of the variance in y can be explained by
 357 the spectral data. This approach introduces realistic variability, reflecting the inherent uncertainties and complexities
 358 often encountered in real-world prediction tasks.

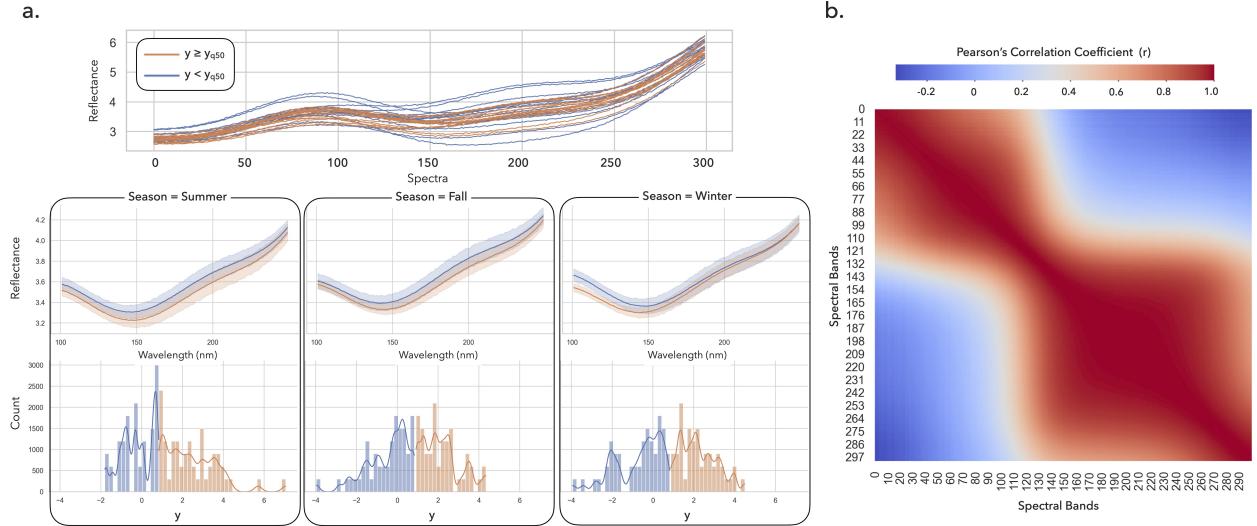


Figure 2: Overview of the simulated spectral dataset. (a) The spectral data matrix X is visualized with the target variable y categorized by its median value. (b) The autocorrelation plot of the spectral data matrix X shows a bi-modular structure with the least pair-wise correlation around the 100th band.

359 The simulated spectral data exhibit a bi-modular autocorrelation structure, with the least pair-wise correlation ($r=0.4$)
 360 observed near the 100th band, which serves as the cutoff between the two modules (Figure 2b). The non-linear
 361 relationship between the spectral data and the target variable y is evident when y is categorized by its median value
 362 (y_{q50}) and visualized in two color groups within the spectral space (Figure 2a). The resulting plot reveals that the
 363 data are not linearly separable in the spectral space, confirming the expected complexity and presenting a challenging
 364 task for model evaluation. Additionally, seasonal effects simultaneously influence both the spectral data and the target
 365 variable. For instance, the spectral reflection is less pronounced around the 150th band during the first season, while
 366 the separability of the two categorical groups decreases around the 250th band in the third season. Furthermore, the
 367 distribution of the target variable varies across seasons: the first season displays a right-skewed distribution, whereas

368 the other two seasons have more symmetric distributions. These seasonal variations introduce additional layers of
 369 complexity, further highlighting the importance of robust evaluation methods for classification models.

370 2.1.3 Real-world spectral dataset

371 This dataset contains spectral data collected across 18 bands ranging from 410 nm to 940 nm, aimed at assessing forage
 372 quality. The spectral data were captured using a SparkFun ESP32 Thing Plus microprocessor paired with a SparkFun
 373 Triad Spectroscopy Sensor (SparkFun Electronics, Niwot, CO). The sensor suite was programmed using Arduino IDE
 374 v2.0.4 (Arduino Core Team, 2024) to export measurements.

375 Forage quality was quantified based on neutral detergent fiber (NDF) content, a critical parameter for evaluating
 376 livestock nutrition. Ground truth NDF values were determined using traditional bench chemistry methods with the
 377 ANKOM 200 fiber analyzer system (ANKOM Technology, Macedon, NY). The dataset comprises 599 samples collected
 378 over three distinct time periods reflecting the seasonal effects on both the spectral data and the NDF response: 189
 379 samples were collected from May to June, 198 from July to August, and 212 from September to October. Sampling
 380 took place weekly between May 1 and October 30, 2023, with two samples collected each week from each of 12 fields.
 381 Within each field, sampling locations were chosen at random and varied from week to week.

382 The spectral data were collected in the field using a handheld sensor, while NDF content was measured in the lab. The
 383 fields were primarily grazed by cattle, with some fields grazed by other species, including sheep and horses. This
 384 dataset provides a comprehensive view of how seasonal variation influences forage quality and spectral characteristics,
 385 offering valuable insights into the dynamics of pasture composition and livestock nutrition.

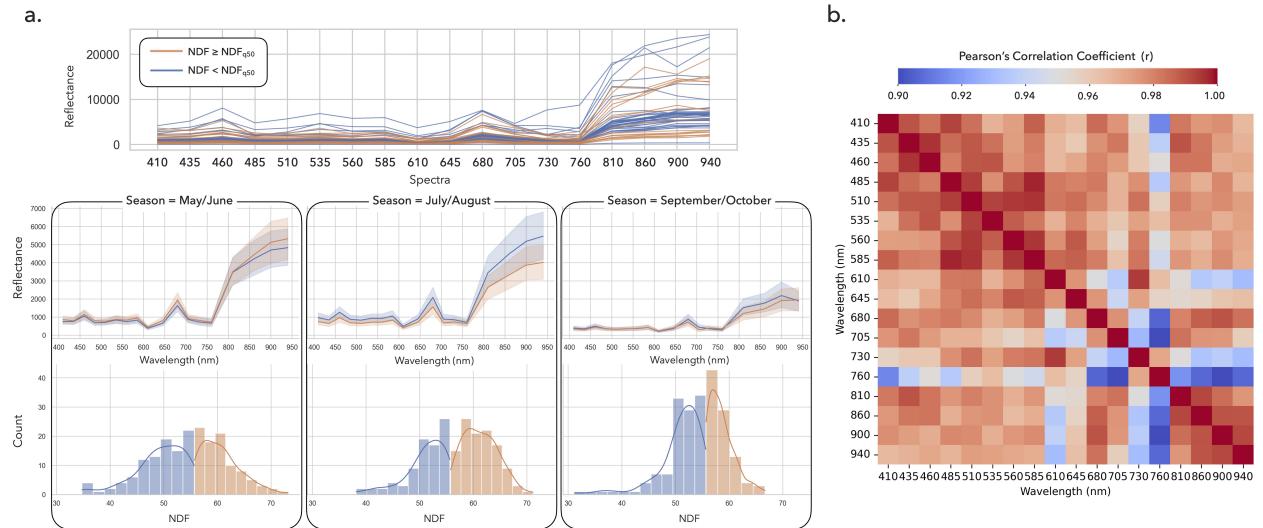


Figure 3: Overview of the real spectral dataset. (a) The spectral data matrix X is visualized with the neutral detergent fiber (NDF) y categorized by its median value. (b) The autocorrelation plot of the spectral data matrix X .

386 A similar examination of the data structure was conducted for both the spectral measurements and the target variable
 387 (Figure 3). The autocorrelation in the spectral measurements is notably stronger than in the simulated dataset, with at

388 least 0.90 for pairwise Pearson correlation coefficients. The seasonal interactions among the spectral measurements are
389 also more pronounced compared to the simulated data. For instance, the spectral reflectance measured in September
390 and October is roughly halved compared to the other seasons. Additionally, a distinct seasonal pattern was observed in
391 the reflectance data beyond 800 nm. In July and August, samples with lower NDF value tend to exhibit higher spectral
392 reflectance, whereas this trend is not evident in May–June or September–October. Moreover, the NDF distribution
393 shows greater variability in July and August, with a higher standard deviation (7.15) compared to May–June (6.19)
394 and September–October (5.23). These observations highlight the stronger seasonal effects and variability in the real
395 dataset compared to the simulated data, providing another example of the challenges in evaluating model performance
396 in real-world agricultural studies.

397 2.2 Study 1: Evaluation bias and variance of cross-validation

398 This study examined the reliability of CV in estimating model performance, with a focus on different performance
399 estimators and their interaction with sample size. It is hypothesized that increasing the number of folds in CV will
400 generally provide a more accurate estimate of model performance but will also lead to increased variance in each
401 estimate, as suggested by the bias-variance trade-off theory. Additionally, sample size is considered a critical factor in
402 reducing the bias difference between estimators, with larger sample sizes expected to mitigate the impact of estimator
403 bias and improve the reliability of performance evaluation.

404 Since K-fold CV employs a fraction (i.e., $K - 1$ folds) of the data for training, it may provide a pessimistic estimate of
405 model performance. Such underestimation is explored in this study by comparing the performance metrics of K-fold
406 CV with K set to 2, 5, and 10, as well as LOOCV where K equals the sample size N, and the "In-Sample" evaluation,
407 which assesses model performance on the same dataset used for training, potentially leading to an overly optimistic
408 bias. To gauge model performance, four metrics are employed: RMSE (Eq. 1.1), MAE (Eq. 1.2), r (Eq. 1.5), and R^2
409 (Eq. 1.6). The evaluation model is a linear regression with ten input features and one output target, all drawn from the
410 null dataset. The sample sizes N are varied among 50, 250, and 500 to explore the dynamics between sample size and
411 performance estimators. Each configuration is repeated across 500 iterations to assess the distribution of evaluation bias
412 and variance.

413 For each iteration, the dataset $\mathcal{D} = (X, Y)$ was sampled as per the simulation's premise. In the case of K-fold CV, the
414 dataset \mathcal{D} was partitioned into K folds in which each fold is $\mathcal{D}_k = (X_k, Y_k)$. For the "In-Sample" approach, partitioning
415 does not occur. The linear model f is trained on the training set \mathcal{D}_{-k} (denoted as $f_{\mathcal{D}_{-k}}$) to estimate regression coefficients
416 β , which then predicts the target variable \hat{Y}_k from the test set \mathcal{D}_k . The procedure of K-fold CV can be expressed as:

$$\begin{aligned}
\text{Training: } Y_{-k} &= f_{\mathcal{D}_{-k}}(X_{-k}) + \epsilon \\
&= X_{-k}\beta + \epsilon \\
\text{Testing: } \hat{Y}_k &= f_{\mathcal{D}_{-k}}(X_k) \\
&= X_k\beta \quad k = 1, 2, \dots, K
\end{aligned} \tag{2.3}$$

417 For the “In-Sample” performance estimator, predictions were made without splitting, as:

$$\begin{aligned}
\text{Training: } Y &= f_{\mathcal{D}}(X) \\
&= X\beta + \epsilon \\
\text{Testing: } \hat{Y} &= f_{\mathcal{D}}(X) \\
&= X\beta
\end{aligned} \tag{2.4}$$

418 Where:

- 419 • X denotes the input regressors sampled from a standard normal distribution $\mathcal{N}(0, 1)$ with dimensions $N \times 10$.
- 420 • Y denotes the target variable sampled from a standard normal distribution $\mathcal{N}(0, 1)$ with dimensions $N \times 1$.
- 421 • X_{-k} and Y_{-k} are the input regressors and target variable in the training set \mathcal{D}_{-k} .
- 422 • X_k denotes the input regressors in the test set \mathcal{D}_k .
- 423 • \hat{Y}_k denotes the predicted target variable in the test set \mathcal{D}_k .
- 424 • β denotes the estimated regression coefficient with dimensions 10×1 .
- 425 • ϵ denotes the error term assumed to be normally distributed.

426 Estimated performance $\mathbb{E}[\hat{g}(f_{\mathcal{D}})]$ was derived by averaging the performance metrics across all K folds as per Eq. S.4.
427 The bias and variance of the evaluation were calculated using Eqs. S.5 and S.6, respectively. To approximate true
428 model performance $G(f_{\mathcal{D}})$, a hundred unseen datasets \mathcal{D}^* were generated identically to \mathcal{D} , and the performance $G(f_{\mathcal{D}})$
429 was estimated by averaging the performance metrics across all \mathcal{D}^* . The detailed steps to compute evaluation bias and
430 variance are provided in the supplementary materials.

431 2.3 Study 2: Model Selection in Cross-Validation

432 The objective of this simulation study is to investigate the impact of improper model selection practices on evaluation
433 bias. Two critical steps in the model selection process are considered: feature selection and hyperparameter tuning.
434 The study hypothesizes that improper model selection — particularly the leakage of test set information during feature
435 selection or hyperparameter tuning — will result in a significant overestimation of model performance.

436 To evaluate this hypothesis, three datasets are utilized: a null dataset with a baseline performance of $r = 0$, a simulated
 437 spectral dataset, and a real spectral dataset. Feature selection is conducted by selecting the top 10 features most strongly
 438 correlated with the target variable, y . The original number of feature candidates varies across datasets, with 1000 for
 439 the null dataset, 300 for the simulated spectral dataset, and 18 for the real spectral dataset.

440 For hyperparameter tuning, the study employs a Support Vector Regression (SVR) model with two hyperparameters:
 441 the kernel function and the regularization parameter (c). The kernel functions considered are linear, sigmoid, and radial
 442 basis function (RBF), while the regularization parameter is set to two values: $c = 1.0$ and $c = 0.01$. The kernel function
 443 determines how the selected features are transformed — either linearly or nonlinearly — to predict the target variable,
 444 y . The regularization parameter c controls the trade-off between minimizing prediction error and model complexity;
 445 a larger c allows for more error but reduces the likelihood of overfitting. In total, six SVR model variants (i.e., three
 446 kernel functions combined with two regularization parameter values) are available for selection during the evaluation
 447 process.

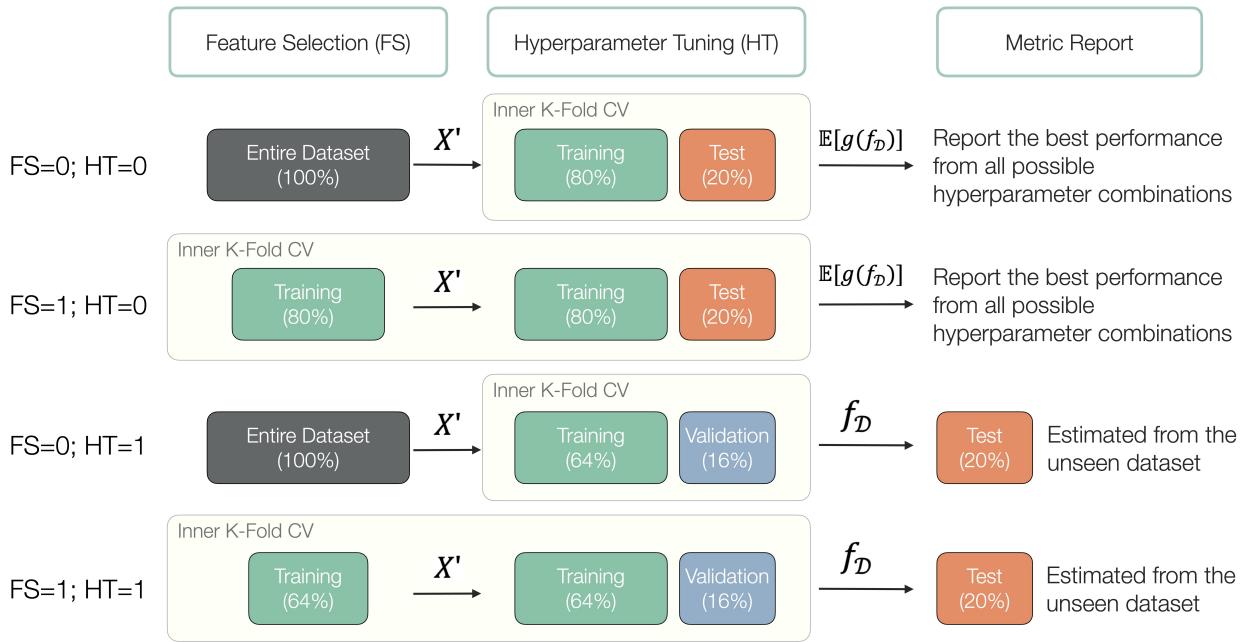


Figure 4: Workflow diagram illustrating four cross-validation strategies of feature selection (FS) and hyperparameter tuning (HT), where 0 denotes incorrect implementation and 1 indicates correct practice. X' is the selected feature subset, $\mathbb{E}[\hat{g}(f_{\mathcal{D}})]$ is the expected generalization performance, $f_{\mathcal{D}}$ is the model trained on the training set without being revealed to the test set.

448 This study introduces notations FS for feature selection and HT for hyperparameter tuning, assigning a binary indicator
 449 (0 or 1) to denote incorrect (0) or correct (1) implementation of model selection. This yields four possible combinations
 450 of model selection strategies: “FS=0; HT=0”, “FS=0; HT=1”, “FS=1; HT=0”, “FS=1; HT=1” (Figure 4). When
 451 FS=0, feature selection precedes cross-validation splitting. If FS=1, feature selection occurs within each fold of the
 452 training set during cross-validation. With hyperparameter tuning, a correct implementation (HT=1) involves splitting
 453 the dataset into training (64%), validation (16%), and test (20%) sets. The model is trained and tuned using the training

and validation sets, respectively, while the test set is reserved for a single evaluation of model performance. Conversely, with $HT=0$, only training (80%) and test (20%) sets are used, risking evaluation bias as the test set informs both training and performance reporting. A 5-fold cross-validation approach was deployed for all strategies. Evaluation bias is measured as the discrepancy between the model selection-influenced performance estimate and the expected generalization performance ($r=0$), using the Pearson correlation coefficient between predicted and observed values. Over 500 sampling iterations, the study assesses the distribution of evaluation bias. A t-test will determine whether the evaluation bias significantly deviates from zero. This experimental setup is designed to quantify the extent of performance overestimation under improper model selection practices and provide insights into its implications for predictive modeling.

2.4 Study 3: Block Effects in Cross-Validation

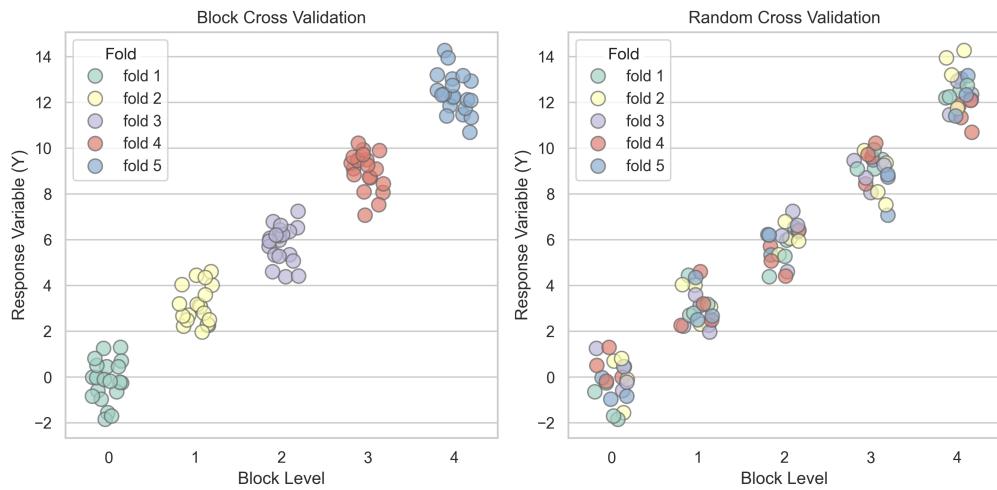


Figure 5: Illustration of fold assignment in block cross-validation (left) and random cross-validation (right). Folds are color-coded, and the block effect is set to 3 in this example.

The objective of this study is to demonstrate how Random CV, which randomly assigns samples to folds without accounting for block effects, can lead to an overestimation of model performance. As a benchmark, the study employs Block CV, where each block is treated as a fold in cross-validation. The hypothesis is that the model performance estimated by Random CV will be significantly higher than that estimated by Block CV.

This study utilizes both simulated and real-world datasets, both of which were collected across multiple seasons, introducing block effects that confound both the predictor features and the response variable. The simulated dataset includes 200 observations per season, distributed equally across seasons, while the real-world dataset also contains approximately 200 observations per season. The block effect in both datasets is defined by the seasonal variation.

The study evaluates two model validation strategies: Block CV and Random CV, both using a 3-fold cross-validation approach. Three folds are used to match the number of seasons in the dataset. In Block CV, each block (i.e., season) is treated as a distinct fold, ensuring that samples from the same block are not split across folds. In Random CV, samples

475 are randomly assigned to folds without consideration of block boundaries (Figure 5). The predictive model used is a
476 random forest regression model, and its performance is assessed using Pearson's correlation coefficient r and RMSE.

477 The simulation is run for 500 iterations, with X (predictor variables) and Y (response variable) resampled in each
478 iteration for the simulation dataset and also the fold assignment to account for variability. A paired t-test is performed
479 to compare the model performance estimates obtained from the two CV strategies, providing a statistical measure of the
480 difference in performance estimation between Random CV and Block CV.

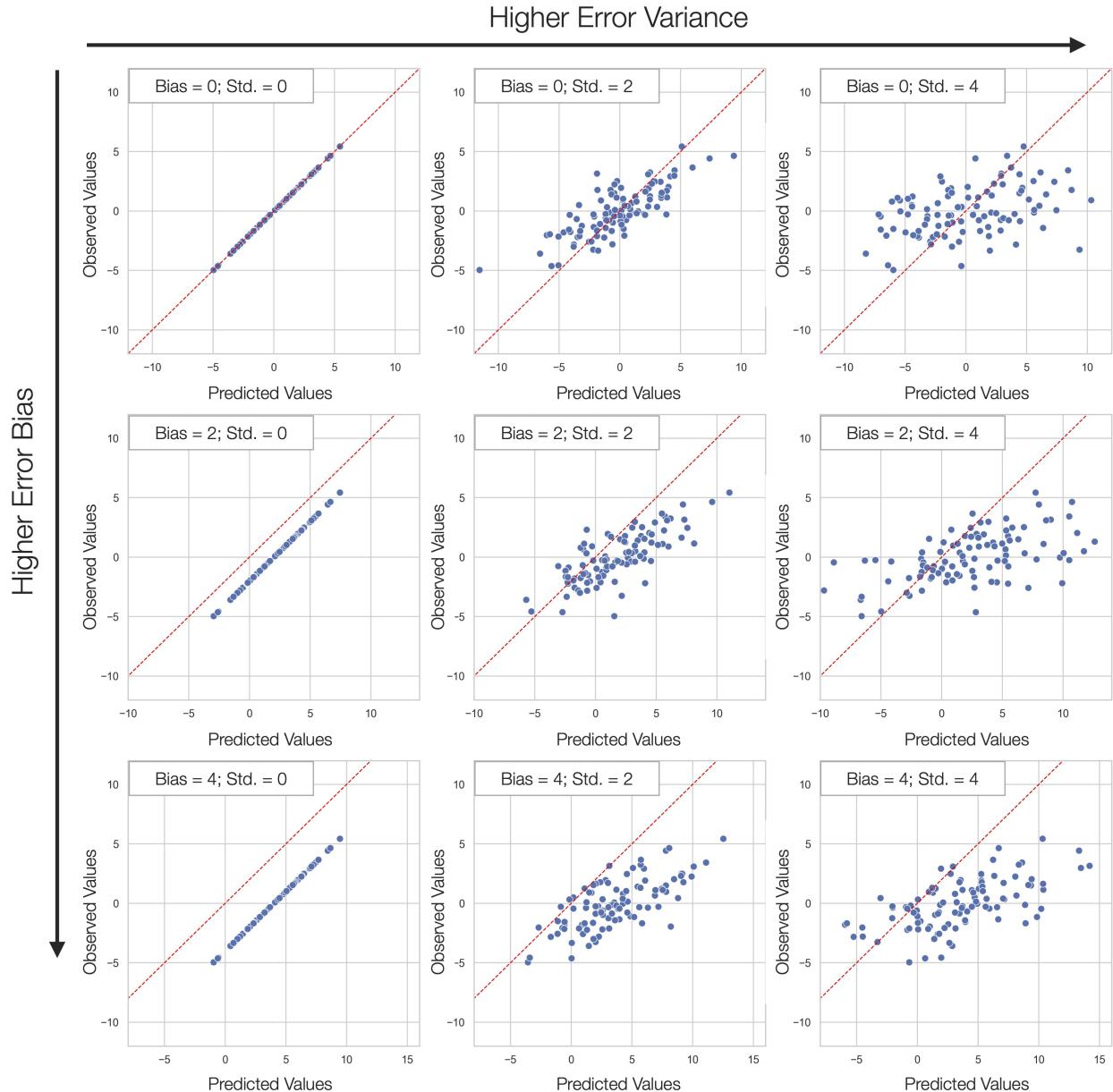
481 **2.5 Study 4: Characteristics of Metrics in Regression Tasks**

Figure 6: Scatter plots illustrating the relationship between predicted and actual values for 9 combinations of bias and variance, with each parameter set to one of three levels: 0, 2, or 4. The red diagonal line represents the ideal prediction line.

482 The objective of this study is to examine how different performance metrics in regression tasks respond to two types
 483 of prediction errors: bias and variance. The study aims to highlight the unique characteristics of each metric, such as
 484 sensitivity to outliers or systematic bias, and provide guidance for selecting appropriate metrics in regression tasks.
 485 Additionally, it seeks to verify the trade-off relationship between bias and variance in prediction errors.

486 To achieve this, six levels of bias and variance are examined: [0, 0.5, 1, 2, 4, 8], forming a total of 36 combinations of
 487 prediction errors. The bias error can be considered as a systematic error that consistently overestimates or underestimates
 488 the ground truth values, while the variance error represents the random fluctuations around the ground truth. Outliers of
 489 prediction errors are considered as extreme cases of variance errors, where the predicted values deviate considerably
 490 from the ground truth. The simulated ground truth values are generated from a normal distribution with a mean of 0 and
 491 a standard deviation of 2, while the predicted values (\hat{y}) are created by adding random errors (ϵ) with specified levels of
 492 bias and variance to the ground truth:

$$\begin{cases} y \sim \mathcal{N}(0, 2) \\ \epsilon \sim \mathcal{N}(b, s) \\ \hat{y} = y + \epsilon \end{cases} \quad (2.5)$$

493 where b represents the bias and s represents the variance of the prediction errors (Figure 6). The choice of a standard
 494 deviation of 2 for the ground truth values is intended to highlight the differences in behavior between the RMSE and
 495 RSR metrics, with RSR being standardized by the standard deviation of the ground truth values while RMSE tracks the
 496 original error scale.

497 The evaluated metrics in this study are categorized into two main groups. Error-based metrics include Root Mean
 498 Squared Error (RMSE), Mean Absolute Error (MAE), and Root Mean Standard Deviation Ratio (RSR). These metrics
 499 focus on quantifying the magnitude of errors in the predictions. Linearity-based metrics, such as Pearson Correlation
 500 Coefficient (r), Coefficient of Determination (R^2), and Concordance Correlation Coefficient (CCC), assess the linear
 501 relationship and agreement between the predicted and actual values. By systematically exploring how each metric
 502 responds to varying levels of bias and variance, this study demonstrates their strengths, limitations, and practical
 503 implications for regression analysis. The findings are intended to guide practitioners in selecting the most appropriate
 504 performance metrics based on their specific modeling objectives and the characteristics of their data.

505 **2.6 Study 5: Characteristics of Metrics in Classification Tasks**

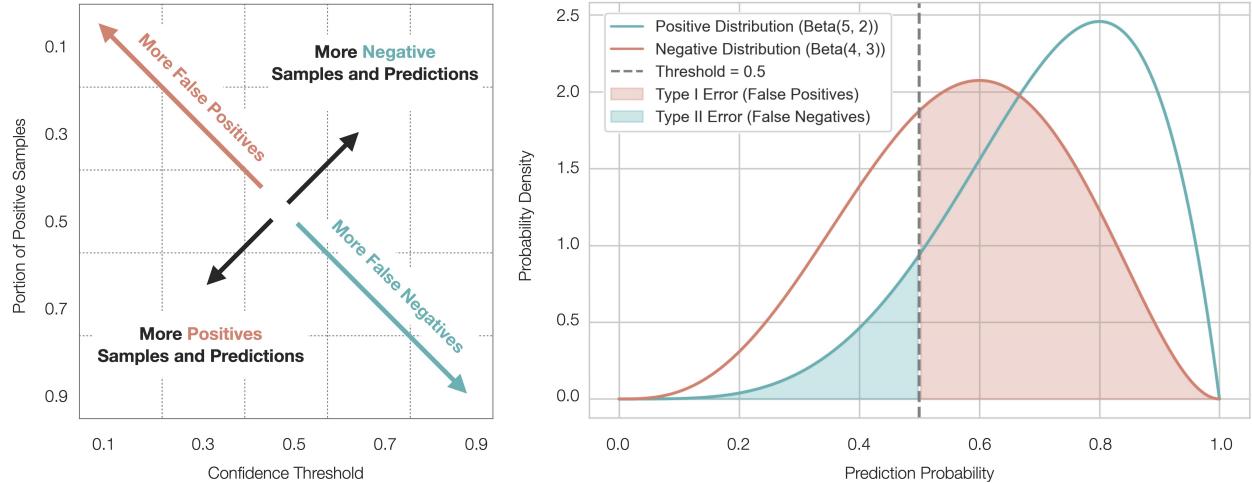


Figure 7: Illustration of the simulation design for evaluating classification metrics under varying balance and confidence thresholds. (a) A 5×5 grid of performance metrics, with each cell representing a unique combination of balance and confidence threshold. (b) The prediction probability distribution for positive and negative samples.

506 The objective of this study is to investigate two critical aspects of evaluating binary classification models: the balance
 507 between positive and negative samples and the choice of confidence threshold. The study aims to explore how these
 508 factors influence the performance metrics used to evaluate classification models and to provide insights into the necessity
 509 of reporting specific metrics together.

510 To achieve this, five levels of balance and five levels of confidence thresholds are examined, forming a total of 25
 511 combinations. The inspected levels are $[0.1, 0.3, 0.5, 0.7, 0.9]$ for both balance and confidence thresholds. The balance
 512 level is determined by the proportion of positive samples, with a balance level of 0.9 indicating that 90% of the samples
 513 are positive. The confidence threshold is used to dichotomize prediction probabilities, where a higher threshold results
 514 in fewer positive predictions by requiring higher certainty for a positive classification. This simulation produces a 5×5
 515 grid of performance metrics, where each cell represents a unique combination of balance and confidence threshold
 516 (Figure 7a). The top-left corner of the grid corresponds to a scenario where positive samples are rare, and the model
 517 uses a low confidence threshold, resulting in a high false positive rate. In contrast, the bottom-right corner represents a
 518 scenario where positive samples are abundant, and the model applies a high confidence threshold, leading to a high
 519 false negative rate. This design contrasts these two extreme cases and hence provides a comprehensive evaluation of
 520 performance metrics across varying conditions.

521 The prediction probability, which represents the likelihood of a given sample being classified as positive by the model,
 522 is simulated using a beta distribution (Figure 7b). For positive samples, the prediction probability is drawn from a beta
 523 distribution with parameters $\alpha = 5$ and $\beta = 2$, resulting in a peak probability around 0.8.

524 For negative samples, the beta distribution has parameters $\alpha = 4$ and $\beta = 3$, with a peak probability around 0.6. This
 525 design creates a scenario where more false positives are expected than false negatives, as the negative samples have high

526 prediction probabilities that overlap with the positive samples. The shaded area in Figure 7b represents the overlap
527 between the two distributions, where the confidence threshold is applied. This overlap introduces two types of errors:
528 the region to the right of the threshold intersecting with the negative distribution represents false positives (Type I error),
529 while the region to the left of the threshold intersecting with the positive distribution represents false negatives (Type II
530 error). This setup highlights the trade-off between these error types based on the choice of the confidence threshold.

531 The metrics evaluated in this study include TPR, TNR, FPR, FNR, sensitivity, specificity, precision, recall, accuracy, F1
532 score, F2 score, and MCC. Although some of these metrics are mathematically equivalent but referred to by different
533 names, this study also highlights the reasons why certain metrics are commonly reported together and discusses their
534 complementary roles in evaluating model performance.

535 **3 Results and Discussion**

536 **3.1 Study 1: The Impact of Estimator Choice and Sample Size on Model Evaluation Reliability**

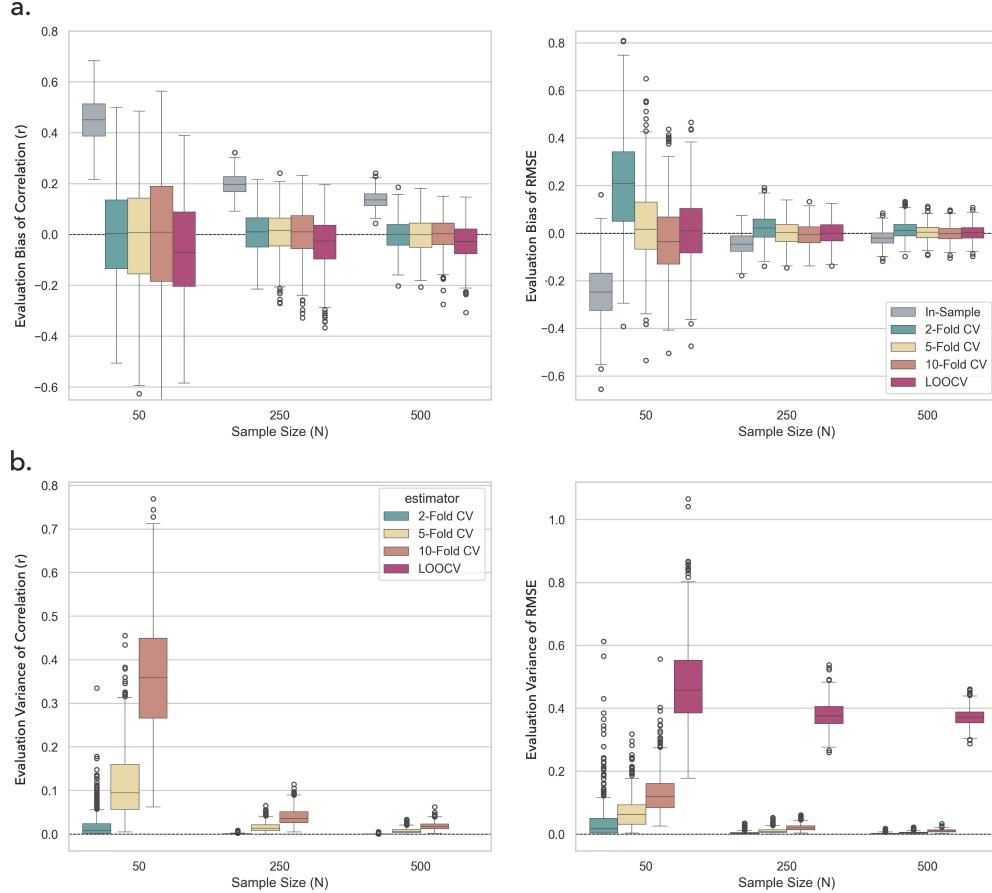


Figure 8: Evaluation bias from 500 sampling iterations on the null dataset with 10 feature variables. Multiple performance estimators across different sample sizes were color-coded. Two metrics: r and RMSE, were displayed in the column facets.

537 The results (Figure 8, Table 4, and Table 5) indicate that both the choice of performance estimator and the sample size
 538 notably influence evaluation reliability, which can be decomposed into bias and variance. Although different numbers of
 539 folds in CV and LOOCV show no substantial differences in bias, they do affect variance. Specifically, as the number of
 540 folds increases, the testing sets become smaller, leading to higher variance. Traditionally, LOOCV has been considered
 541 an unbiased estimator for error-based metrics such as R^2 , RMSE, and MAE. In this study, LOOCV generally follows
 542 that expectation, except for a few cases at certain sample sizes. Interestingly, when the ratio of sample size to number of
 543 features is sufficiently high (e.g., 25 when $N = 250$), other K -fold estimators can deliver comparable accuracy and
 544 bias, offering a more computationally efficient alternative to LOOCV.

545 However, a key finding emerges with correlation-based metrics (e.g., r): LOOCV tends to underestimate model
 546 performance and exhibits a pessimistic bias. At $N = 250$ and $N = 500$, LOOCV's bias on r can be 10 to 30 times

547 larger than other K -fold estimators. In contrast, in-sample (or apparent) estimation, while conventionally deemed the
 548 most biased due to information leakage from the testing set, can surprisingly achieve comparable reliability at larger
 549 sample sizes. For instance, at $N = 250$, the bias of in-sample estimation is only 0.099 for R^2 and -0.044 for RMSE,
 550 which is less biased than all K -fold CV estimators for R^2 and less biased than 2-fold CV for RMSE with a smaller
 551 sample size of 50.

Table 4: Evaluation bias (mean \pm std) for the metrics from 500 sampling iterations. The minimum bias given the same sample size is highlighted in bold.

| Metric | Estimator | N=50 | N=250 | N=500 |
|--------|------------|------------------------------------|------------------------------------|------------------------------------|
| r | In-Sample | 0.449 \pm 0.088 | 0.198 \pm 0.043 | 0.137 \pm 0.033 |
| | 2-Fold CV | 0.004\pm0.184 | 0.009 \pm 0.082 | -0.001 \pm 0.061 |
| | 5-Fold CV | -0.012 \pm 0.209 | 0.006 \pm 0.088 | -0.001 \pm 0.067 |
| | 10-Fold CV | -0.011 \pm 0.254 | 0.003\pm0.094 | 0.000\pm0.065 |
| | LOOCV | -0.070 \pm 0.203 | -0.035 \pm 0.098 | -0.031 \pm 0.071 |
| R^2 | In-Sample | 0.515 \pm 0.207 | 0.099 \pm 0.037 | 0.053 \pm 0.020 |
| | 2-Fold CV | -0.694 \pm 0.642 | -0.044 \pm 0.071 | -0.017 \pm 0.034 |
| | 5-Fold CV | -0.401 \pm 0.409 | -0.024 \pm 0.049 | -0.007\pm0.026 |
| | 10-Fold CV | -0.940 \pm 0.857 | -0.046 \pm 0.052 | -0.014 \pm 0.024 |
| | LOOCV | -0.013\pm0.256 | 0.009\pm0.039 | 0.008 \pm 0.020 |
| RMSE | In-Sample | -0.244 \pm 0.116 | -0.044 \pm 0.044 | -0.020 \pm 0.032 |
| | 2-Fold CV | 0.215 \pm 0.226 | 0.022 \pm 0.056 | 0.013 \pm 0.036 |
| | 5-Fold CV | 0.035 \pm 0.158 | 0.002 \pm 0.047 | 0.004 \pm 0.033 |
| | 10-Fold CV | -0.024 \pm 0.149 | -0.006 \pm 0.046 | -0.001\pm0.033 |
| | LOOCV | 0.012\pm0.144 | 0.001\pm0.046 | 0.003 \pm 0.033 |
| MAE | In-Sample | -0.195 \pm 0.096 | -0.037 \pm 0.037 | -0.017 \pm 0.028 |
| | 2-Fold CV | 0.180 \pm 0.190 | 0.017 \pm 0.047 | 0.010 \pm 0.030 |
| | 5-Fold CV | 0.049 \pm 0.134 | 0.004 \pm 0.039 | 0.003 \pm 0.028 |
| | 10-Fold CV | 0.022 \pm 0.127 | 0.002 \pm 0.038 | 0.002 \pm 0.028 |
| | LOOCV | 0.011\pm0.119 | -0.001\pm0.038 | 0.001\pm0.028 |

Table 5: Evaluation variance (mean \pm std) for the metrics from 500 sampling iterations. The minimum variance given the same sample size is highlighted in bold.

| Metric | Estimator | N=50 | N=250 | N=500 |
|--------|------------|-----------------------------------|-----------------------------------|-----------------------------------|
| r | 2-Fold CV | 0.019\pm0.030 | 0.001\pm0.001 | 0.000\pm0.001 |
| | 5-Fold CV | 0.117 \pm 0.081 | 0.016 \pm 0.011 | 0.008 \pm 0.005 |
| | 10-Fold CV | 0.362 \pm 0.131 | 0.040 \pm 0.019 | 0.019 \pm 0.009 |
| R^2 | 2-Fold CV | 0.859 \pm 2.876 | 0.003\pm0.006 | 0.001\pm0.001 |
| | 5-Fold CV | 0.743\pm1.391 | 0.008 \pm 0.009 | 0.002 \pm 0.002 |
| | 10-Fold CV | 7.164 \pm 37.486 | 0.018 \pm 0.019 | 0.003 \pm 0.002 |
| RMSE | 2-Fold CV | 0.041\pm0.068 | 0.004\pm0.005 | 0.002\pm0.003 |
| | 5-Fold CV | 0.070 \pm 0.050 | 0.010 \pm 0.008 | 0.005 \pm 0.004 |
| | 10-Fold CV | 0.130 \pm 0.067 | 0.021 \pm 0.010 | 0.010 \pm 0.005 |
| | LOOCV | 0.477 \pm 0.131 | 0.379 \pm 0.041 | 0.371 \pm 0.029 |
| MAE | 2-Fold CV | 0.030\pm0.053 | 0.003\pm0.003 | 0.001\pm0.002 |
| | 5-Fold CV | 0.052 \pm 0.039 | 0.008 \pm 0.005 | 0.004 \pm 0.003 |
| | 10-Fold CV | 0.100 \pm 0.052 | 0.015 \pm 0.007 | 0.008 \pm 0.004 |
| | LOOCV | 0.477 \pm 0.131 | 0.379 \pm 0.041 | 0.371 \pm 0.029 |

552 Further examination confirms that a higher number of folds in CV generally reduces bias for error-based metrics (RMSE,
 553 MAE) because training sets become more representative of the total data. Yet, correlation-based metrics can display
 554 divergent trends under the same conditions. In LOOCV, evaluating a single data point at a time makes its variance
 555 particularly evident for RMSE, as single-point predictions inherently fluctuate more. Consequently, across all sample
 556 sizes tested, LOOCV consistently exhibits greater variance than lower-fold CV (e.g., 2-fold, 5-fold). Nonetheless, bias
 557 and variance across all estimators converge as sample size grows (e.g., $N = 500$).

558 In conclusion, performance estimation reliability depends strongly on the interplay between the estimation method, the
 559 metric in use, and the sample size. Larger sample sizes typically reduce both bias and variance, thereby improving the
 560 trustworthiness of model evaluations. While LOOCV often provides less biased estimates for error-based metrics, it can
 561 severely underestimate correlation-based metrics and suffers from higher variance. K -fold CV methods present a more
 562 computationally manageable solution for large datasets and can match LOOCV's performance when the sample size
 563 is sufficiently large relative to the number of features. Ultimately, selecting the most appropriate evaluation strategy
 564 should be based on practical considerations—such as available sample size, computational resources, and the specific
 565 metrics of interest—to ensure robust and reliable model assessments.

566 **3.2 Study 2: Misuse of Model Selection Can Lead to Over-Optimistic Performance Estimates**

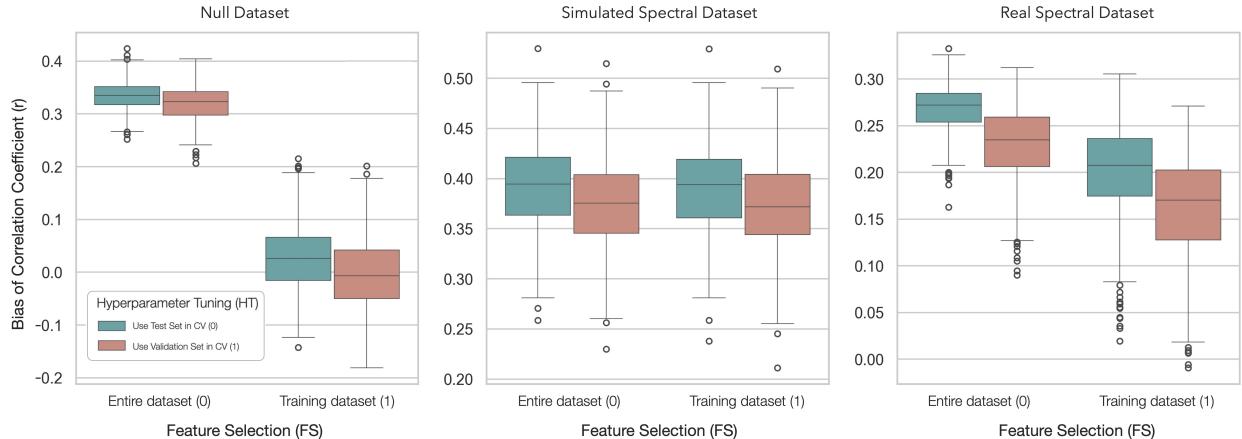


Figure 9: The evaluation bias of the four model selection strategies.

567 With the exception of the feature selection procedure in the simulated spectral dataset, all other procedures that
 568 erroneously incorporate the testing set into the model selection process substantially inflate evaluation bias (Figure 9
 569 and Table 6). However, the magnitude of this inflation varies across datasets. In the null dataset, incorrectly performing
 570 feature selection leads to roughly a 30% performance inflation; on the real dataset, a similar practice results in a 6%
 571 inflation. By contrast, the simulated spectral dataset shows no measurable inflation, with only a negligible 0.02% change
 572 in r . This minimal effect may stem from the degree to which the selected features contribute to predictive accuracy. If a
 573 model's performance is comparable to that obtained through random feature selection, then leveraging the testing set
 574 for feature selection does not necessarily boost performance. Such a scenario is more plausible in datasets with high

575 multicollinearity, where multiple features correlate strongly, making any subset of features effectively representative of
 576 the entire feature space.

Table 6: ANOVA results of how each procedure affects the evaluation bias measured in the correlation coefficient (r). FS: Feature Selection, HT: Hyperparameter Tuning. DF: Degree of Freedom, SS: Sum of Squares, MS: Mean Squares. Significant p-values (< 0.05) are highlighted in bold.

(a) Dataset: Null dataset, Metric: r

| Factor | DF | SS | MS | F-value | p-value |
|----------|------|-------|-------|----------|------------------|
| FS | 1 | 49.72 | 49.72 | 20073.41 | < 1e-6 |
| HT | 1 | 0.24 | 0.24 | 97.83 | < 1e-6 |
| FS:HT | 1 | 0.03 | 0.03 | 14.33 | < 1e-6 |
| Residual | 1996 | 4.94 | 0.00 | — | — |

(b) Dataset: Simulated spectral dataset, Metric: r

| Factor | DF | SS | MS | F-value | p-value |
|----------|------|----------|----------|---------|------------------|
| FS | 1 | 1.87e-03 | 1.87e-03 | 1.03 | 0.391 |
| HT | 1 | 1.64e-01 | 1.64e-01 | 91.10 | < 1e-6 |
| FS:HT | 1 | 2.85e-08 | 2.85e-08 | 0.00 | 0.997 |
| Residual | 1996 | 3.60e+00 | 1.80e-03 | — | — |

(c) Dataset: Real spectral dataset, Metric: r

| Factor | DF | SS | MS | F-value | p-value |
|----------|------|------|------|---------|------------------|
| FS | 1 | 2.31 | 2.31 | 1198.87 | < 1e-6 |
| HT | 1 | 0.73 | 0.73 | 382.76 | < 1e-6 |
| FS:HT | 1 | 0.00 | 0.00 | 0.20 | 0.648 |
| Residual | 1996 | 3.85 | 0.00 | — | — |

577 On the other hand, using the entire dataset to perform hyperparameter tuning substantially inflates performance in all
 578 three datasets—by approximately 2.5% in the null dataset, 1.5% in the simulated spectral dataset, and nearly 4% in
 579 the real spectral dataset. It is worth noting that these estimates arise from a relatively small search space of only six
 580 hyperparameter combinations. In contemporary machine learning, hyperparameter spaces can be far larger, especially
 581 for deep learning models in which the architectures themselves are highly configurable, involving potentially millions
 582 of parameters. Even minor alterations (e.g., changing the kernel size in a convolutional layer) may markedly affect
 583 model performance in such complex settings [40].

584 Collectively, these findings underscore the importance of rigorous CV practices in model selection, particularly for
 585 feature selection and hyperparameter tuning, to achieve accurate performance estimates and robust generalizability in
 586 predictive modeling.

587 **3.3 Study 3: Overlooking Experimental Block Effects Can Lead to Biased Model Performance Estimates**

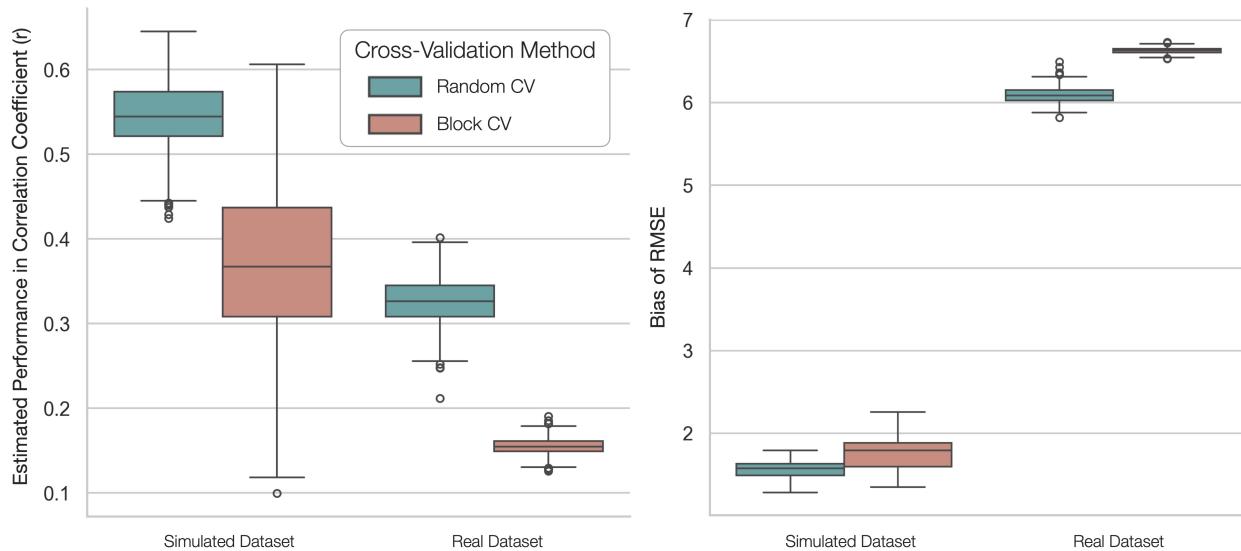


Figure 10: Bias in model performance estimation by Block CV and Random CV across 1000 iterations. The dashed line represents the null hypothesis that the mean performance estimate is zero.

588 Performance inflation is evident in both the simulated and real spectral datasets, which inherently exhibit seasonal
 589 variation as block effects (Figure 10 and Table 7). Ignoring this seasonal variation leads to a notable overestimation
 590 of model performance, as reflected by a 17.5% bias in r for the simulated dataset and 17.1% for the real dataset. A
 591 similar pattern emerges for RMSE, with a 15.5% bias in the simulated dataset and an 11.1% bias in the real dataset.
 592 The ANOVA results further support these findings, since all four tests show significant differences (p -value < 0.001)
 593 between the two methods on the estimated model performance.
 594 These observations highlight the importance of closely examining any identifiable sources of variation in the experiment
 595 and aligning evaluation strategies with the model's intended real-world application. Variation such as the seasonal block
 596 effects demonstrated here can simultaneously influence both the predictive features and the response variable. If the
 597 model is intended for deployment in a new block, such as a future season for which no information is available prior to
 598 deployment, using block CV is essential to ensure generalizability. On the other hand, if the model is designed for a
 599 closed environment where all possible blocks are already represented, random CV offers a more efficient strategy.

Table 7: ANOVA results for different datasets and metrics. DF: Degree of Freedom, SS: Sum of Squares, MS: Mean Squares. Significant p-values (< 0.05) are highlighted in bold.

(a) Dataset: Simulated spectral dataset, Metric: r

| Factor | DF | SS | MS | F-value | p-value |
|----------|-----|------|------|---------|------------------|
| method | 1 | 9.61 | 9.61 | 2122.69 | < 1e-6 |
| Residual | 998 | 4.52 | 0.00 | — | — |

(b) Dataset: Real spectral dataset, Metric: r

| Factor | DF | SS | MS | F-value | p-value |
|----------|-----|------|------|----------|------------------|
| method | 1 | 8.64 | 8.64 | 29744.48 | < 1e-6 |
| Residual | 998 | 0.29 | 0.00 | — | — |

(c) Dataset: Simulated spectral dataset, Metric: $RMSE$

| Factor | DF | SS | MS | F-value | p-value |
|----------|-----|-------|-------|---------|------------------|
| method | 1 | 11.57 | 11.57 | 559.59 | < 1e-6 |
| Residual | 998 | 20.64 | 0.02 | — | — |

(d) Dataset: Real spectral dataset, Metric: $RMSE$

| Factor | DF | SS | MS | F-value | p-value |
|----------|-----|-------|-------|----------|------------------|
| method | 1 | 88.40 | 88.40 | 26768.87 | < 1e-6 |
| Residual | 998 | 3.29 | 0.00 | — | — |

600 3.4 Study 4: Characteristics of Metrics in Regression Tasks

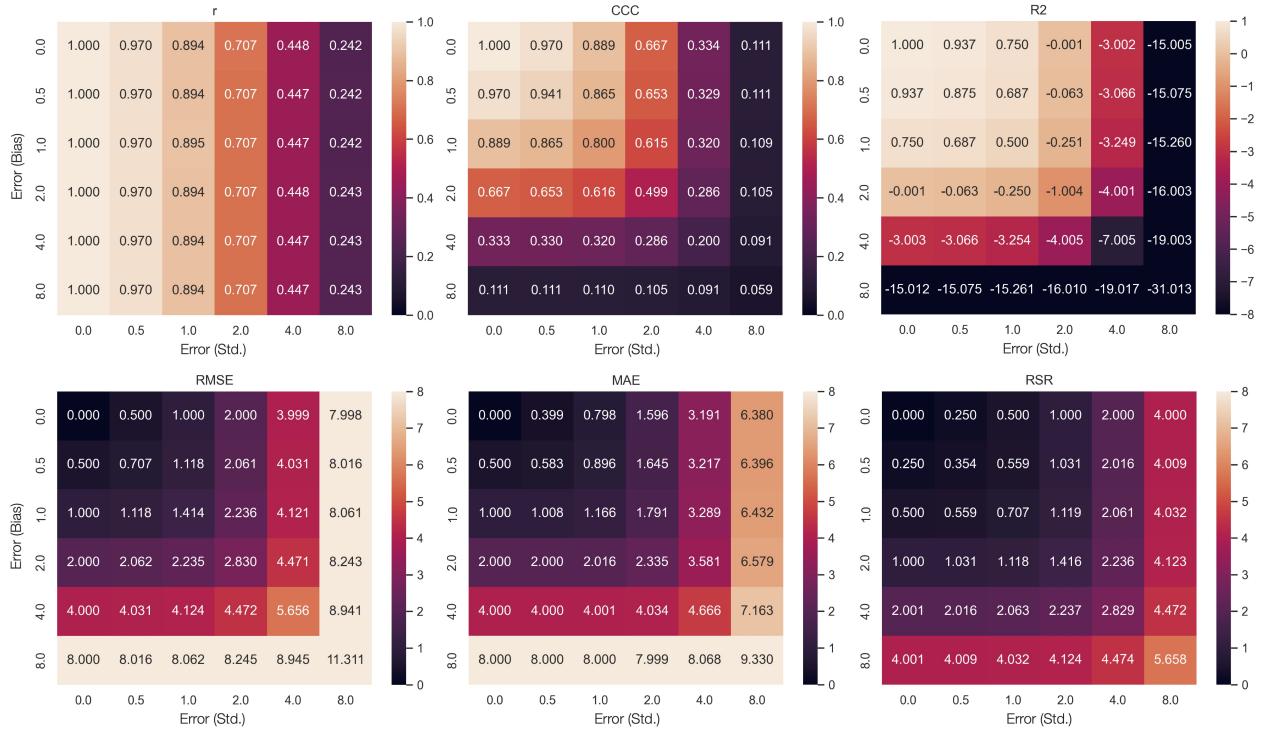


Figure 11: The heatmap of different metrics in response to the bias and variance error in regression tasks.

601 Metrics in regression tasks show distinct trend in response to different combination of error bias and variance (Figure
 602 11). Except for r and MAE, all metrics show symmetric response in the grid matrix, supporting the trade-off relationship
 603 between error bias and variance. For example, R^2 show the same 0.75 values at position of (1, 4) and (4, 1), indicating
 604 that the model with 0 bias and 2 variance is equivalent to the model with 2 bias and 0 variance. Both r and CCC are
 605 linearity-based metric measuring the correlation between predicted and actual values but present different trend. Since
 606 r has been standardized by the standard deviation of both predicted and actual values, it is completely invariant to the
 607 bias error that could reflect the scale of the prediction error. Whereas CCC is sensitive to both the bias and variance,
 608 hence providing a multi-faced evaluation of the correlation-based evaluation. In this study, CCC also show a better
 609 interpretability than r . According to the result, when $CCC > 0.5$, the total squared error (i.e., the sum of squared
 610 bias and variance) is no larger than $4 + 4 = 8$, which is twice the ground truth variance ($2^2 = 4$). Furthermore, when
 611 $CCC > 0.8$, the error now must be less than the ground truth variance. Compared to the ambiguous interpretation of r ,
 612 these benchmarks provide a straightforward guideline to translate the correlation concept into the prediction error scale.
 613 R^2 is a popular metric in machine learning community, and it also provides a good interpretability. When $R^2 = 0$, it
 614 suggested the model is no better than the mean prediction. The mean prediction, in this heatmap, can be understood as
 615 the prediction has same levels of error as the ground truth, which is also found in the position of (4, 1) and (1, 4) in the
 616 heatmap. Furthermore, when $R^2 = -1$, it suggested the prediction is worse than the mean prediction by a unit of the
 617 ground truth variance, which is also found in the position of (4, 4) where both the error bias and variance are 2 and

618 make the total squared error 8, which is twice the ground truth variance. Although share the similar interpretability with
619 CCC, when encountering outlier prediction, R^2 inflated rapidly from -3 to -15, for example, when the bias/variance error
620 changed from 4 to 8. The same error change only result in 0.22 decrease in CCC, which is always in a range between 0
621 to 1 and making it a more stable and comparable metric across different prediction context.

622 The error-based metrics, RMSE and MAE are usually compared for their robustness to outliers. With the square error
623 term in RMSE, it is more sensitive to the large prediction error. Interestingly, when comparing the metric values in
624 responding to the increase of bias error, RMSE and MAE remain exact the same trend: 1 unit of bias error increase
625 will also result in 1 unit of metric value increase. The different responsiveness only observed along the variance error
626 axis, where MAE only increase 0.8 unit while RMSE increase 1 unit per variance error increase by 1. This result
627 provide a detailed dissection of the robustness of these two metrics to the prediction error, and suggest that MAE is more
628 tolerant to the error variance than RMSE, making it robust to the outlier prediction. RSR is a metric that standardize
629 the RMSE by the standard deviation of the ground truth values. In this study, the ground truth value were generated
630 through a normal distribution with a standard deviation of 2. This setting make the RSR metric to be a direct comparison
631 with RMSE: RMSE reflect the original error scale, while RSR standardize the scale to the ground truth variance,
632 making its value always halved of the RMSE value. This nature make RSR a good metric when comparing the model
633 performance across multiple datasets with different variance level. RSR provide a uniform standard to evaluate the
634 model performance, but it also keep the capability of tracking the error amplitude that linearity-based metrics like r and
635 CCC lack.

636 **3.5 Study 5: Characteristics of Metrics in Classification Tasks**

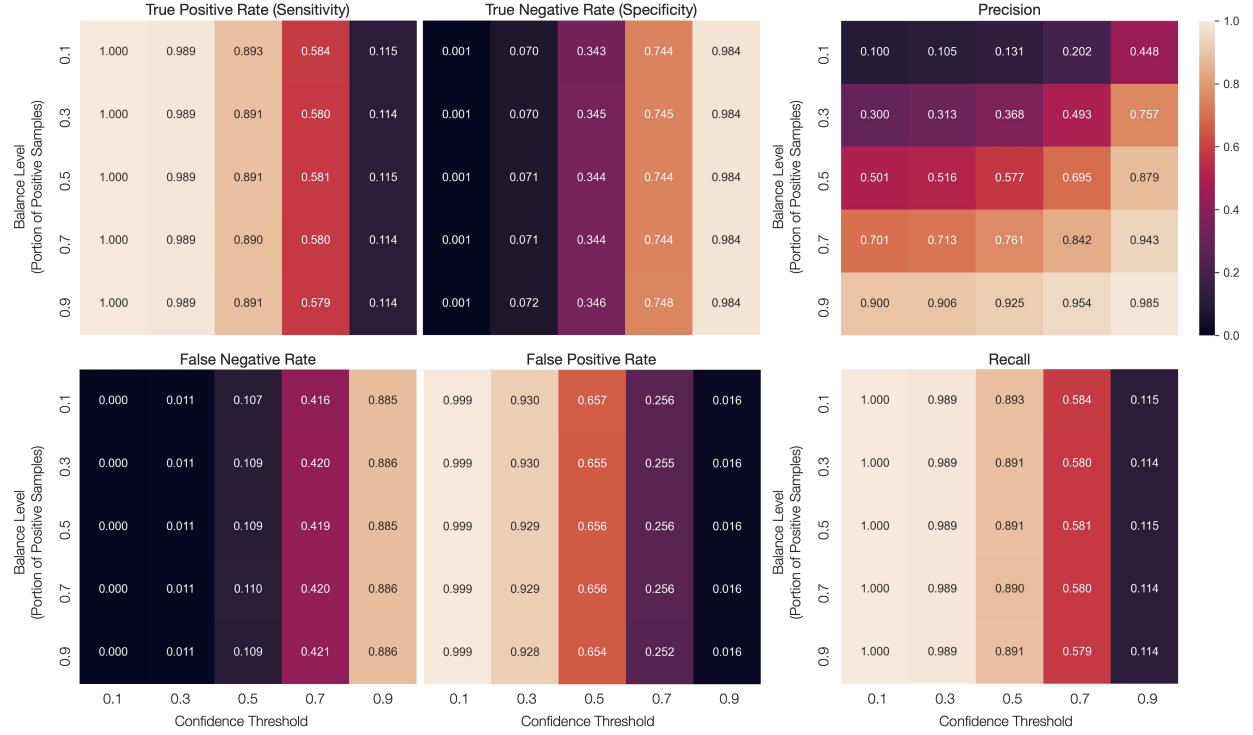


Figure 12: A 5×5 grid of performance metrics, with each cell representing a unique combination of balance and confidence threshold. The inspected metrics focus on one sample distribution at a time.

- 637 Except for precision, metrics that focus on a single sample distribution (i.e., either actual positives or actual negatives)
 638 remain invariant to changes in class imbalance (Figure 12). For instance, regardless of the proportion of positive samples
 639 in the dataset, the true positive rate (TPR) remains at 0.989 when the confidence threshold is 0.3. This stability arises
 640 because these metrics consider only one sample distribution at a time, rendering them unaffected by shifts in the overall
 641 balance between positive and negative samples.
- 642 Given the nature, multiple metrics are often used in combination, as each captures different aspects of model performance
 643 that may be overlooked by a single metric. In this study, where false positives are generally more prevalent than false
 644 negatives, reporting TPR alone does not sufficiently reflect the model’s capability, given that it focuses solely on the
 645 actual positive class. Including the true negative rate (TNR) is therefore important to account for the contribution of
 646 false positives among actual negative samples. When examining both TPR and TNR, the overall lower performance of
 647 TNR across various confidence thresholds helps to highlight the model’s tendency to produce more false positives than
 648 false negatives. Notably, the false positive rate (FPR) complements TNR (since $FPR + TNR = 1$), and the false negative
 649 rate (FNR) complements TPR (since $FNR + TPR = 1$). These pairs can thus be used interchangeably, offering flexibility
 650 in how the model’s classification errors are reported and interpreted.
- 651 Precision as the only metric that considers the predicted positive samples is sensitive to both class imbalance and
 652 the confidence threshold. When using precision to report performance, the balance level should be considered as the

baseline for interpreting the results. For example, when balance = 0.7 where 70% of samples are positive, a naive model predicting all samples as positive would achieve a precision of 0.7, which is higher than the precision of 0.5 achieved by a model that predicts samples at random. From this heatmap, the precision values in the first column shows almost no difference from the designated balance level, indicating that a model with a low confidence threshold is likely to predict most samples as positive and has no actual predictive power. In contrast, with balance=0.1 and confident threshold=0.9, a 0.448 precision should be considered a good performance, as the baseline precision is 0.1 and the extra 0.348 precision is gained by correctly predicting positive samples from a negative-majority dataset.

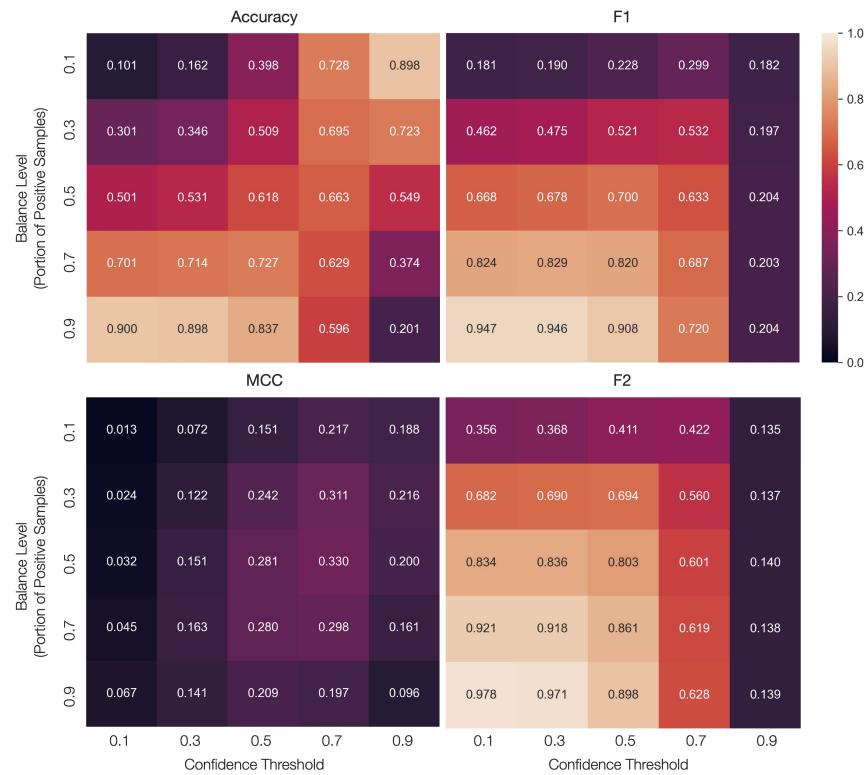


Figure 13: A 5×5 grid of performance metrics, with each cell representing a unique combination of balance and confidence threshold. The inspected metrics focus on multiple sample distributions simultaneously.

Metrics that capture multiple facets of model performance provide deeper insights (Figure 13). In the heatmap used in this study, the four corners represent extreme scenarios. The top-left corner corresponds to a dataset dominated by negative samples but with predominantly positive predictions (thus yielding many false positives), while the bottom-right corner reflects a dataset dominated by positive samples but with predominantly negative predictions (thus yielding many false negatives). The remaining two corners illustrate cases in which both predictions and actual labels are overwhelmingly of the same class, conditions under which a naive model that predicts only the majority class could still achieve high performance. Consequently, an effective metric must capture a model's true predictive power and move beyond these forms of "background performance."

668 Accuracy provides an overall assessment of model performance and exhibits lower values in the top-left corner, where
669 false positives are more pronounced. This observation aligns with the study's setting, in which false positives are more
670 frequent than false negatives. Although the top-right and bottom-left corners both exhibit high accuracy, these outcomes
671 reflect "background performance" because the accuracy is essentially the same as the proportion of the majority class.
672 This limitation arises because accuracy does not account for other facets of model performance. For example, in the
673 top-right corner, where 90% of the samples are negative, the model's high confidence threshold causes it to predict all
674 samples as negative, yielding 0.90 accuracy. However, this result captures only the true negative rate and overlooks the
675 fact that the true positive rate is very low.

676 The F1 score addresses this shortcoming by considering both the true positive rate (recall) and precision, providing a
677 more balanced measure of performance. In this study for the same top-right corner, the F1 score is 0.182, reflecting
678 sensitivity to both false positives and false negatives. A variant of the F1 score, known as F2, is also presented for
679 comparison. By definition, F2 places greater emphasis on true positive rate, making it more sensitive to false negatives.
680 This property is evident in the simulation, where F2 outperforms F1 in the top-left corner, where false negatives have
681 a stronger influence on overall performance. Such variations of the F1 score can be selected based on the specific
682 priorities of the application context.

683 Finally, the MCC is often considered the most balanced metric, making it more stringent in evaluating model perfor-
684 mance. Any substantial deficiency, such as a high false negative rate, is reflected in a lower MCC. Although focusing
685 solely on MCC does not indicate precisely which aspect of the model requires improvement, it conveniently highlights
686 where the model performs optimally. In this study, MCC exceeds 0.30 only in a narrow region where the dataset is
687 balanced (balance =0.5) and the threshold is moderately high (0.7), pointing to an ideal operational setting for this
688 model.

689 4 Conclusion

690 In conclusion, this study presents a comprehensive evaluation of five simulation experiments, uncovering critical
691 insights into the interplay of performance estimators, metrics, and contextual factors that influence model evaluation
692 reliability. The findings highlight the nuanced impact of estimator choices and sample sizes on bias and variance,
693 emphasizing that while traditional estimators such as LOOCV can be reliable for error-based metrics, they may
694 severely underestimate correlation-based metrics under certain conditions. The misuse of model selection processes
695 was shown to substantially inflate performance estimates, reinforcing the importance of adhering to rigorous cross-
696 validation practices. Additionally, the role of experimental block effects in biasing performance estimates underscores
697 the necessity of aligning evaluation strategies with real-world applications. In both regression and classification
698 tasks, metric characteristics vary significantly, with different metrics offering complementary perspectives on model
699 performance. For regression tasks, CCC and R^2 provided interpretable benchmarks for understanding prediction errors,
700 while MAE demonstrated greater robustness to variance compared to RMSE. In classification tasks, metrics such as
701 precision, F1 score, and MCC were shown to capture different facets of performance, with MCC offering a balanced,
702 stringent evaluation across multiple dimensions. Collectively, these findings stress the importance of tailoring model
703 evaluation strategies to specific research and application contexts to ensure robust, reliable, and interpretable results.

704 5 Acknowledgement

705 The author James Chen expresses his gratitude to Drs. Zhiwu Zhang, Hao Cheng, Gota Morota, and Gonzalo Ferreira
706 for their insightful discussions that partially contributed to this study. The authors declare no conflicts of interest.

707 **References**

- 708 [1] Hao Cheng, Dorian J. Garrick, and Rohan L. Fernando. Efficient strategies for leave-one-out cross validation for
709 genomic best linear unbiased prediction. *Journal of Animal Science and Biotechnology*, 8(1):38, May 2017.
- 710 [2] I. D. E. van Dixhoorn, R. M. de Mol, J. T. N. van der Werf, S. van Mourik, and C. G. van Reenen. Indicators of
711 resilience during the transition period in dairy cows: A case study. *Journal of Dairy Science*, 101(11):10271–10282,
712 November 2018.
- 713 [3] T. Hastie, R. Tibshirani, and J.H. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and
714 Prediction*. Springer series in statistics. Springer, 2009.
- 715 [4] Gavin C. Cawley and Nicola L.C. Talbot. On Over-fitting in Model Selection and Subsequent Selection Bias in
716 Performance Evaluation. *The Journal of Machine Learning Research*, 11:2079–2107, August 2010.
- 717 [5] Arthur E. Hoerl and Robert W. Kennard. Ridge Regression: Biased Estimation for Nonorthogonal Problems.
718 *Technometrics*, 12(1):55–67, 1970. Publisher: [Taylor & Francis, Ltd., American Statistical Association, American
719 Society for Quality].
- 720 [6] Robert Tibshirani. Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society:
721 Series B (Methodological)*, 58(1):267–288, January 1996.
- 722 [7] Harris Drucker, Chris J. C. Burges, Linda Kaufman, Alex Smola, and Vladimir Vapnik. Support vector regression
723 machines. In *Proceedings of the 9th International Conference on Neural Information Processing Systems*, NIPS’96,
724 pages 155–161, Cambridge, MA, USA, December 1996. MIT Press.
- 725 [8] Hervé Abdi. Partial Least Square Regression PLS-Regression. *Encyclopedia of social sciences research methods*,
726 pages 792–795, 2003.
- 727 [9] Leo Breiman. Random Forests. *Machine Learning*, 45(1):5–32, October 2001.
- 728 [10] Yann LeCun. Generalization and Network Design strategies. 1989.
- 729 [11] Morteza H. Ghaffari, Amirhossein Jahanbekam, Hassan Sadri, Katharina Schuh, Georg Dusel, Cornelia Prehn,
730 Jerzy Adamski, Christian Koch, and Helga Sauerwein. Metabolomics meets machine learning: Longitudinal
731 metabolite profiling in serum of normal versus overconditioned cows and pathway analysis. *Journal of Dairy
732 Science*, 102(12):11561–11585, December 2019.
- 733 [12] G. Rovere, G. de los Campos, A. L. Lock, L. Worden, A. I. Vazquez, K. Lee, and R. J. Tempelman. Prediction of
734 fatty acid composition using milk spectral data and its associations with various mid-infrared spectral regions in
735 Michigan Holsteins. *Journal of Dairy Science*, 104(10):11242–11258, October 2021.
- 736 [13] C. A. Becker, A. Aghalari, M. Marufuzzaman, and A. E. Stone. Predicting dairy cattle heat stress using machine
737 learning techniques. *Journal of Dairy Science*, 104(1):501–524, January 2021.

- 738 [14] B. Lahart, S. McParland, E. Kennedy, T.M. Boland, T. Condon, M. Williams, N. Galvin, B. McCarthy, and
739 F. Buckley. Predicting the dry matter intake of grazing dairy cows using infrared reflectance spectroscopy analysis.
740 *Journal of Dairy Science*, 102(10):8907–8918, October 2019.
- 741 [15] Tiago Bresolin and João R. R. Dórea. Infrared Spectrometry as a High-Throughput Phenotyping Technology to
742 Predict Complex Traits in Livestock Systems. *Frontiers in Genetics*, 11, 2020.
- 743 [16] C. Grelet, E. Froidmont, L. Foldager, M. Salavati, M. Hostens, C. P. Ferris, K. L. Ingvartsen, M. A. Crowe, M. T.
744 Sorensen, J. A. Fernandez Pierna, A. Vanlierde, N. Gengler, and F. Dehareng. Potential of milk mid-infrared
745 spectra to predict nitrogen use efficiency of individual dairy cows in early lactation. *Journal of Dairy Science*,
746 103(5):4435–4445, May 2020.
- 747 [17] I. Adriaens, N. C. Friggins, W. Ouweltjes, H. Scott, B. Aernouts, and J. Statham. Productive life span and
748 resilience rank can be predicted from on-farm first-parity sensor time series but not using a common equation
749 across farms. *Journal of Dairy Science*, 103(8):7155–7171, August 2020.
- 750 [18] Lucio F. M. Mota, Diana Giannuzzi, Vittoria Bisutti, Sara Pegolo, Erminio Trevisi, Stefano Schiavon, Luigi Gallo,
751 David Fineboym, Gil Katz, and Alessio Cecchinato. Real-time milk analysis integrated with stacking ensemble
752 learning as a tool for the daily prediction of cheese-making traits in Holstein cattle. *Journal of Dairy Science*,
753 105(5):4237–4255, May 2022.
- 754 [19] Roii Spoliansky, Yael Edan, Yisrael Parmet, and Ilan Halachmi. Development of automatic body condition scoring
755 using a low-cost 3-dimensional Kinect camera. *Journal of Dairy Science*, 99(9):7714–7725, September 2016.
- 756 [20] Sun Yukun, Huo Pengju, Wang Yujie, Cui Ziqi, Li Yang, Dai Baisheng, Li Runze, and Zhang Yonggen. Automatic
757 monitoring system for individual dairy cows based on a deep learning framework that provides identification via
758 body parts and estimation of body condition score. *Journal of Dairy Science*, 102(11):10140–10151, November
759 2019.
- 760 [21] X. Song, E.A.M. Bokkers, P.P.J. Van Der Tol, P.W.G. Groot Koerkamp, and S. Van Mourik. Automated body
761 weight prediction of dairy cows using 3-dimensional vision. *Journal of Dairy Science*, 101(5):4448–4459, May
762 2018.
- 763 [22] C. Xavier, Y. Le Cozler, L. Depuille, A. Caillot, A. Lebreton, C. Allain, J. M. Delouard, L. Delattre, T. Luginbuhl,
764 P. Faverdin, and A. Fischer. The use of 3-dimensional imaging of Holstein cows to estimate body weight and
765 monitor the composition of body weight change throughout lactation. *Journal of Dairy Science*, 105(5):4508–4519,
766 May 2022.
- 767 [23] P. Mäntysaari, E.A. Mäntysaari, T. Kokkonen, T. Mehtio, S. Kajava, C. Grelet, P. Lidauer, and M.H. Lidauer.
768 Body and milk traits as indicators of dairy cow energy status in early lactation. *Journal of Dairy Science*,
769 102(9):7904–7916, September 2019.

- 770 [24] M. Frizzarin, I. C. Gormley, D. P. Berry, T. B. Murphy, A. Casa, A. Lynch, and S. McParland. Predicting cow
771 milk quality traits from routinely available milk spectra using statistical machine learning methods. *Journal of*
772 *Dairy Science*, 104(7):7438–7447, July 2021.
- 773 [25] J. A. D. R. N. Appuhamy, J. V. Judy, E. Kebreab, and P. J. Kononoff. Prediction of drinking water intake by dairy
774 cows. *Journal of Dairy Science*, 99(9):7191–7205, September 2016.
- 775 [26] R. A. de Souza, R. J. Tempelman, M. S. Allen, W. P. Weiss, J. K. Bernard, and M. J. VandeHaar. Predicting
776 nutrient digestibility in high-producing dairy cows. *Journal of Dairy Science*, 101(2):1123–1135, February 2018.
- 777 [27] J. R. R. Dórea, G. J. M. Rosa, K. A. Weld, and L. E. Armentano. Mining data from milk infrared spectroscopy to
778 improve feed intake predictions in lactating dairy cows. *Journal of Dairy Science*, 101(7):5878–5889, July 2018.
- 779 [28] L. I. Lin. A concordance correlation coefficient to evaluate reproducibility. *Biometrics*, 45(1):255–268, March
780 1989.
- 781 [29] Edward J. Jones, Thomas F. A. Bishop, Brendan P. Malone, Patrick J. Hulme, Brett M. Whelan, and Patrick
782 Filippi. Identifying causes of crop yield variability with interpretive machine learning. *Computers and Electronics*
783 *in Agriculture*, 192:106632, January 2022.
- 784 [30] N. W. O’Leary, D. T. Byrne, A. H. O’Connor, and L. Shalloo. Invited review: Cattle lameness detection with
785 accelerometers. *Journal of Dairy Science*, 103(5):3895–3911, May 2020.
- 786 [31] J. Stojkov, G. Bowers, M. Draper, T. Duffield, P. Duivenvoorden, M. Groleau, D. Haupstein, R. Peters,
787 J. Pritchard, C. Radom, N. Sillett, W. Skippon, H. Trépanier, and D. Fraser. Hot topic: Management of cull
788 dairy cows—Consensus of an expert consultation in Canada. *Journal of Dairy Science*, 101(12):11170–11174,
789 December 2018.
- 790 [32] Maher Alsaad, Mahmoud Fadul, and Adrian Steiner. Automatic lameness detection in cattle. *The Veterinary*
791 *Journal*, 246:35–44, April 2019.
- 792 [33] X. Kang, X. D. Zhang, and G. Liu. Accurate detection of lameness in dairy cattle with computer vision: A new
793 and individualized detection strategy based on the analysis of the supporting phase. *Journal of Dairy Science*,
794 103(11):10628–10638, November 2020.
- 795 [34] S. J. Denholm, W. Brand, A. P. Mitchell, A. T. Wells, T. Krzyzelewski, S. L. Smith, E. Wall, and M. P. Coffey.
796 Predicting bovine tuberculosis status of dairy cows from mid-infrared spectral data of milk using deep learning.
797 *Journal of Dairy Science*, 103(10):9355–9367, October 2020.
- 798 [35] S.A. Kandeel, A.A. Megahed, M.H. Ebeid, and P.D. Constable. Ability of milk pH to predict subclinical mastitis
799 and intramammary infection in quarters from lactating dairy cattle. *Journal of Dairy Science*, 102(2):1417–1427,
800 February 2019.
- 801 [36] Davide Chicco and Giuseppe Jurman. The advantages of the Matthews correlation coefficient (MCC) over F1
802 score and accuracy in binary classification evaluation. *BMC Genomics*, 21(1):6, January 2020.

- 803 [37] J. M. Bowen, M. J. Haskell, G. A. Miller, C. S. Mason, D. J. Bell, and C-A. Duthie. Early prediction of
804 respiratory disease in preweaning dairy calves using feeding and activity behaviors. *Journal of Dairy Science*,
805 104(11):12009–12018, November 2021.
- 806 [38] Maxime Metz, Alessandra Biancolillo, Matthieu Lesnoff, and Jean-Michel Roger. A note on spectral data
807 simulation. *Chemometrics and Intelligent Laboratory Systems*, 200:103979, May 2020.
- 808 [39] Chun-Peng James Chen, Yang Hu, Xianran Li, Craig F. Morris, Stephen Delwiche, Arron H. Carter,
809 Camille Steber, and Zhiwu Zhang. An independent validation reveals the potential to predict Hag-
810 berg–Perten falling number using spectrometers. *The Plant Phenome Journal*, 6(1):e20070, 2023. _eprint:
811 <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ppj2.20070>.
- 812 [40] Barret Zoph and Quoc V. Le. Neural Architecture Search with Reinforcement Learning, February 2017.
813 arXiv:1611.01578 [cs].

814 **Appendix**815 **Cross Validation**

816 Model cross validation aims to evaluate how well a given model generalizes to an independent dataset that it has not
 817 seen during the training process. The most common method is K-fold cross-validation (**K-fold CV**). To implement the
 818 K-fold CV, the available dataset, denoted as \mathcal{D} , is partitioned into K equally sized folds. We can express the dataset as
 819 below:

$$\begin{aligned}\mathcal{D} &= \{(X, Y)\} \\ &= \{(X_1, Y_1), (X_2, Y_2), \dots, (X_K, Y_K)\}\end{aligned}\tag{S.1}$$

820 where $X \in \mathbb{R}^{n \times p}$ represents the input features, and $Y \in \mathbb{R}^{n \times 1}$ symbolizes the ground truth labels for a single target
 821 variable. The value of n corresponds to the total number of samples, while p represents the number of features. In
 822 each iteration of the K-fold CV, a single fold is reserved as the test set, $\mathcal{D}_{\text{test}}$ (or \mathcal{D}_k), to act as unseen data, while the
 823 remaining folds make up the training set $\mathcal{D}_{\text{train}}$ (or \mathcal{D}_{-k}):

$$\begin{aligned}\mathcal{D}_{\text{train}} &= \mathcal{D}_{-k} \\ &= \{(X_1, Y_1), (X_2, Y_2), \dots, (X_{k-1}, Y_{k-1}), (X_{k+1}, Y_{k+1}), \dots, (X_K, Y_K)\} \\ \mathcal{D}_{\text{test}} &= \mathcal{D}_k \\ &= \{(X_k, Y_k)\}\end{aligned}\tag{S.2}$$

824 After splitting the dataset into \mathcal{D}_{-k} and \mathcal{D}_k , the examined model f is trained on the training set \mathcal{D}_{-k} and denoted as $f_{\mathcal{D}_{-k}}$.
 825 The hold-out test set \mathcal{D}_k is then used to evaluate the model performance $\hat{g}(f_{\mathcal{D}_{-k}})$, which is defined by comparing the
 826 predicted labels $\hat{Y}_k = f_{\mathcal{D}_{-k}}(X_k)$ with the true labels Y_k using a performance metric \mathcal{L} (e.g., RMSE or R^2):

$$\begin{aligned}\hat{g}(f_{\mathcal{D}_{-k}}) &= \mathcal{L}(Y_k, \hat{Y}_k) \\ &= \mathcal{L}(Y_k, f_{\mathcal{D}_{-k}}(X_k))\end{aligned}\tag{S.3}$$

827 To estimate the generalization performance of a model $\mathbb{E}[\hat{g}(f_{\mathcal{D}})]$, the K-fold CV procedure is repeated K times until
 828 each fold has been used as the test set \mathcal{D}_k once. The entire dataset \mathcal{D} is leveraged to calculate the average prediction
 829 performance over all K folds. The model's generalization performance can be expressed as:

$$\begin{aligned}\mathbb{E}[\hat{g}(f_{\mathcal{D}})] &= \mathbb{E}[\hat{g}(f_{\mathcal{D}_k})] \\ &= \frac{1}{K} \sum_{k=1}^K \hat{g}(f_{\mathcal{D}_k})\end{aligned}\tag{S.4}$$

830 It is noted that $\mathbb{E}[\hat{g}(f_{\mathcal{D}})]$ is equivalent to $\mathbb{E}[\hat{g}(f_{\mathcal{D}_k})]$ in K-fold CV. It is because the $\mathbb{E}[\hat{g}(f_{\mathcal{D}})]$ is estimated by averaging
 831 all $\hat{g}(f_{\mathcal{D}_k})$ over K folds, which is also the definition of $\mathbb{E}[\hat{g}(f_{\mathcal{D}_k})]$.

832 Cross Validation Bias and Variance

833 The true generalization performance of the model $G(f_{\mathcal{D}})$ can only be approximated by averaging the performance
 834 metrics over infinite unseen datasets. However, in practice, the dataset \mathcal{D} is finite and therefore, there is always a bias
 835 when using a finite dataset to estimate $G(f_{\mathcal{D}})$. The bias is known as validation bias:

$$\text{Bias} = \mathbb{E}[\hat{g}(f_{\mathcal{D}})] - G(f_{\mathcal{D}})\tag{S.5}$$

836 For example, if RMSE is used as the performance metric, a positive validation bias suggests that the model validation
 837 procedure concludes a pessimistic estimation of the model performance, since the true performance is expected to be
 838 lower than the estimated performance. Another aspect of model validation is the variance of the estimated performance.
 839 For example, in a 5-fold cross-validation, there are five estimates of the model performance. The variance among these
 840 five estimates is known as validation variance. A high validation variance suggests that the performance is sensitive to
 841 the choice of the test set \mathcal{D}_k , which may be caused by a small sample size or an over-complex model. The validation
 842 variance can be defined as:

$$\begin{aligned}\text{Variance} &= \mathbb{E}[(\hat{g}(f_{\mathcal{D}_k}) - \mathbb{E}[\hat{g}(f_{\mathcal{D}})])^2] \\ &= \mathbb{E}[\hat{g}^2(f_{\mathcal{D}_k}) - 2\hat{g}(f_{\mathcal{D}_k})\mathbb{E}[\hat{g}(f_{\mathcal{D}})] + \mathbb{E}^2[\hat{g}(f_{\mathcal{D}})]] \\ &= \mathbb{E}[\hat{g}^2(f_{\mathcal{D}_k})] - 2\mathbb{E}[\hat{g}(f_{\mathcal{D}_k})]\mathbb{E}[\hat{g}(f_{\mathcal{D}})] + \mathbb{E}^2[\hat{g}(f_{\mathcal{D}})] \\ &= \mathbb{E}[\hat{g}^2(f_{\mathcal{D}_k})] - \mathbb{E}^2[\hat{g}(f_{\mathcal{D}})]\end{aligned}\tag{S.6}$$

843 Combining the Equations S.5 and S.6, the mean squared error (MSE) of the model validation can be decomposed as:

$$\begin{aligned}
\text{MSE} &= \mathbb{E}[(\hat{g}(f_{D_k}) - G(f_D))^2] \\
&= \mathbb{E}[\hat{g}^2(f_{D_k})] - 2\mathbb{E}[\hat{g}(f_{D_k})]G(f_D) + G^2(f_D) + \\
&\quad \mathbb{E}^2[\hat{g}(f_{D_k})] - \mathbb{E}^2[\hat{g}(f_{D_k})] \\
&= (\mathbb{E}^2[\hat{g}(f_{D_k})] - 2\mathbb{E}[\hat{g}(f_{D_k})]G(f_D) + G^2(f_D)) + \\
&\quad (\mathbb{E}[\hat{g}^2(f_{D_k})] - \mathbb{E}^2[\hat{g}(f_{D_k})]) \\
&= (\mathbb{E}[\hat{g}(f_{D_k})] - G(f_D))^2 + (\mathbb{E}[\hat{g}^2(f_{D_k})] - \mathbb{E}^2[\hat{g}(f_{D_k})]) \\
&= (\mathbb{E}[\hat{g}(f_D)] - G(f_D))^2 + (\mathbb{E}[\hat{g}^2(f_D)] - \mathbb{E}^2[\hat{g}(f_D)]) \\
&= \text{Bias}^2 + \text{Variance}
\end{aligned} \tag{S.7}$$

844 **Hyperparameter**

845 Here are the loss functions for ordinary least squares (OLS), ridge regression, and LASSO regression, respectively:

$$\mathcal{L}_{\text{OLS}}(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 \tag{S.8}$$

$$\mathcal{L}_{\text{ridge}}(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \tag{S.9}$$

$$\mathcal{L}_{\text{LASSO}}(\beta) = \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \sum_{j=1}^p |\beta_j| \tag{S.10}$$

846 Where x_i and y_i represent the i th row of the design matrix X and the response vector Y , respectively. The term n
 847 denotes the sample size, and β is the coefficient vector. All three models aim to find the optimal β that minimizes their
 848 respective loss function, \mathcal{L} . In the regularized models (i.e., ridge and LASSO regression), the vector length of β is
 849 penalized in the loss function.

850 **Squared Correlation Coefficient r^2 and Determination Coefficient R^2**

851 The squared Pearson correlation coefficient, r^2 , is not necessarily equivalent to the coefficient of determination, R^2 .
 852 This equivalence holds true specifically in the context of least squares regression when the same model and data are
 853 used for both fitting and evaluation. However, this may not be the case when the model is assessed using new data.
 854 To demonstrate the equivalence between r^2 and R^2 under these specific conditions, we begin by assuming that the
 855 covariance between the predicted values \hat{Y} and the residuals ϵ is zero:

$$\begin{aligned}
\text{cov}(Y, \hat{Y}) &= \text{cov}(\hat{Y} + \epsilon, \hat{Y}) \\
&= \text{cov}(\hat{Y}, \hat{Y}) + \text{cov}(\hat{Y}, \epsilon) \\
&= \text{var}(\hat{Y}) + \text{cov}(\hat{Y}, \epsilon) \\
&= \text{var}(\hat{Y})
\end{aligned} \tag{S.11}$$

856 With the assumption that $\bar{\hat{Y}} = \bar{Y}$, which typically holds when $\mathbb{E}[\epsilon] = 0$, the squared correlation coefficient r^2 is
 857 expressed as follows:

$$\begin{aligned}
r^2 &= \frac{\text{cov}^2(Y, \hat{Y})}{\text{var}(Y)\text{var}(\hat{Y})} \\
&= \frac{\text{var}(\hat{Y})^2}{\text{var}(Y)\text{var}(\hat{Y})} \\
&= \frac{\text{var}(\hat{Y})}{\text{var}(Y)} \\
&= \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{\hat{Y}})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\
&= \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{\hat{Y}})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\
&= \frac{SS_{\text{regression}}}{SS_{\text{total}}} \\
&= R^2
\end{aligned} \tag{S.12}$$

858 where $SS_{\text{regression}}$ is the variation explained by the model and SS_{total} is the total sum of squares. Each Y_i and \hat{Y}_i are the
 859 i th elements of the actual response vector Y and the predicted response vector \hat{Y} , while \bar{Y} and $\bar{\hat{Y}}$ are their respective
 860 means. This proof highlights that under certain assumptions, r^2 and R^2 can indeed be equivalent, but such conditions
 861 are specific to least squares regression where errors are normally distributed and predictions are unbiased estimates of
 862 the actual values.

863 **Study 1: Raw evaluation bias and variance**