# Overview of leakage scenarios in supervised machine learning

**Short title: Leakage in ML**

**Authors:**

Sasse, L.[1,a,b,c]; Nicolaisen-Sobesky, E.[1,a,b]; Dukart, J.[a,b]; Eickhoff, S.B.[a,b]; Götz, M.[d,e]; Hamdan, S.[a,b]; Komeyer, V.[a,b,f]; Kulkarni, A.[g]; Lahnakoski, J.[a,b]; Love, B.C.[h,i,j]; Raimondo, F.[a,b]; Patil, K.R*.[a,b,k].

a.  Institute of Neuroscience and Medicine, Brain and Behaviour (INM-7), Forschungszentrum Jülich, Jülich, Germany

b.  Institute of Systems Neuroscience, Medical Faculty, Heinrich-Heine-University Düsseldorf, Düsseldorf, Germany

c.  Max Planck School of Cognition, Stephanstrasse 1a, Leipzig, Germany

d.  Division of Experimental Radiology, Department for Diagnostic and Interventional Radiology, University Hospital Ulm, Ulm, Germany

e.  Experimental Radiology, University Ulm, Ulm, Germany

f.  Department of Biology, Faculty of Mathematics and Natural Sciences, Heinrich Heine University Duesseldorf, Duesseldorf, Germany

g.  Principal Global Services, Pune, India

h.  Department of Experimental Psychology, University College London, London, UK

i.  The Alan Turing Institute, London, UK

j.  European Lab for Learning & Intelligent Systems (ELLIS)

k.  Koita Centre for Digital Health, IIT Bombay, Mumbai 400076, India

*Corresponding author: Kaustubh R Patil <k.patil@fz-juelich.de

[1]These authors contributed equally

**Abstract:**

29  Machine learning (ML) provides powerful tools for predictive modeling. ML's popularity stems

30  from the promise of sample-level prediction with applications across a variety of fields from

31  physics and marketing to healthcare. However, if not properly implemented and evaluated,

32  ML pipelines may contain leakage typically resulting in overoptimistic performance estimates

33  and failure to generalize to new data. This can have severe negative financial and societal

34  implications. Our aim is to expand understanding associated with causes leading to leakage

35  when designing, implementing, and evaluating ML pipelines. Illustrated by concrete examples,

36  we provide a comprehensive overview and discussion of various types of leakage that may

37  arise in ML pipelines.

**List of abbreviations:**

AUROC:  Area under the receiver operating characteristic curve

CV: Cross validation

DOME: Data, optimization, model and evaluation

FC: Functional connectivity

HCP-YA: Human connectome project - young adult

I.I.D.: Independently and identically distributed

KNN: K-nearest neighbors

49     LOO: Leave-one-out

50     MAE: Mean absolute error

51     ML: Machine learning

52     PCA: Principal component analysis

53     PLS: Partial least squares

54     SVM: support vector machine

55

## 1. Introduction

Machine learning (ML) has become a popular approach to make predictions, aid decision making, and gain insights into complex data in numerous scientific fields. Various methodologies like supervised, unsupervised, generative, and reinforcement learning define the ML landscape, each with its own unique strengths and applications. In supervised learning, the machine learns a function that links input to output by utilizing labeled training data where the correct output is known, derived from example input-output pairs (1). Unsupervised learning deals with unlabeled data. The machine must figure out the correct answer without being told about a ground truth and must therefore discover patterns and structures in the input data (e.g. using clustering) (1). Generative learning is a machine learning approach centered on generating novel data samples. This technique is commonly applied in tasks like producing images, texts, and various other data types (2). Reinforcement learning involves an agent interacting with an environment, taking actions, and receiving rewards or penalties. Through repeated interactions, the model autonomously learns the optimal strategy to maximize rewards, relying less on external guidance for output determination (3).

However, despite the due use and applicability of those methods, supervised learning remains prominent for predictive modeling with applications in various domains including health-care, physics, and climate science (4–12). This is not only because supervised learning is well suited to learn from tabular data ubiquitously found in scientific domains, but also because easy-to-use software libraries with hundreds of learning algorithms and data wrangling tools have lowered the entry barrier for supervised ML-based analyses (e.g. scikit-learn (13) and tidymodels (14)). These collaborative advancements in accessible tools, expanding datasets, and evolving methodologies demonstrate the considerable promise of supervised ML applications to drive transformative innovation across diverse problem domains. This paper, therefore, is concerned with supervised ML.

81    Despite the availability of easy-to-use ML software, most applications still require assembling

82    a custom ML-based data analysis pipeline satisfying unique considerations in terms of data

83    preprocessing, feature engineering, (hyper)parameter tuning, and model selection. While end-

84    to-end tools exist, opting for these often sacrifices control for convenience (e.g. (15)).

85    Therefore, implementing a correct ML pipeline and drawing valid conclusions from the ensuing

86    results remains challenging, and prone to errors. This challenge extends beyond technical

87    aspects, impacting the interpretability and trustworthiness of the outcomes. Handcrafted

88    pipelines, although demanding, afford practitioners the precision, control and insight required

89    for complex data analysis scenarios. Additionally, the evolving nature of data and algorithmic

90    advancements continually reshapes best practices, necessitating a balance between

91    automation and custom solutions to ensure accuracy and relevance in analyses. Striking this

92    balance remains pivotal for robust, reliable, and impactful ML applications.

93    ML models are powerful, and they are adept at exploiting any available information. Thus, it

94    falls on the practitioner to ensure that the modeling approach is reliable and valid. As we will

95    discuss, even simple ML pipelines, if not properly implemented and interpreted, can lead to

96    drastically wrong interpretations and severely problematic conclusions. These issues extend

97    far beyond academic debates; they hold immense societal relevance. Widespread adoption

98    of flawed practices in machine learning can exact substantial societal and economic costs,

99    underscoring the urgency to rectify and mitigate these risks (16). Similar to the replication

100   crisis that recently engulfed the statistics communities and much of the applied sciences,

101   owing to misunderstanding and –intentional or unintentional– misuse of $p$-values from null

102   hypothesis significance testing (17,18), misunderstandings and malpractice in ML can lead to

103   its own replication crisis (16,19) with severe negative financial and societal ramifications (20).

104   It must be noted that reproducibility of a ML pipeline is not sufficient to resolve this, as a

105   reproducible ML pipeline could be still incorrect in inference. Addressing such ML pitfalls is

106   essential to improve the quality and trustworthiness of ML-based data analyses, and

107   consequently will lead to better applications and foster societal acceptance. While previous

108    works have addressed several pitfalls in ML-based analysis (7,21–27), only a few have

109    covered the wide range of threats posed by leakage (16,28–30) (also see John Langford:

110    https://hunch.net/?p=22).

111    Data leakage is one of the most common and most critical types of error when applying ML.

112    Data leakage refers to the leakage of "illegitimate" information into the training process of a

113    ML model (16,28). For example, leakage occurs when the model gets to learn from information

114    about the supposed unseen test set. Therefore, a fair evaluation of the generalization error is

115    not possible, as the test set does not really represent new, unseen data anymore. This likely

116    means that any estimate of the error will be overly optimistic (16). The threat of data leakage

117    can be exemplified by the case of a recent study that claimed high accuracy (91%) in predicting

118    suicidality in youth using neuroimaging data (31). Such a model would be of high clinical

119    relevance and could provide valuable insights about underlying brain phenotypes. However,

120    this paper was retracted because it relied on leakage-prone feature selection leading to an

121    overfitted model and erroneous interpretations (20,32). Hence, the threat posed by leakage in

122    ML pipelines severely affects realistic estimation of generalization performance, insights

123    gained, and deployment.

124    Data leakage is a widespread pitfall on ML pipelines across numerous scientific fields (16).

125    This highlights the importance of raising awareness among a very broad community of

126    researchers encompassing different fields. Recent studies have contributed to raising

127    awareness on the threats of data leakage (16,28–30). However, data leakage is a complex

128    issue that can happen in numerous ways, and often in subtle ways which are difficult to

129    pinpoint. Moreover, even though some types of leakage are recognized and well-discussed in

130    the literature, such as illegitimate use of the targets of the test data (28), many scenarios of

131    data leakage remain unexplored. All these subtleties and under exploration of data leakage

132    make its detection a complicated and tricky task. The limited awareness of its threats and its

133    widespread presence in many fields, underscores the urgency to raise awareness regarding

134    a wide array of data leakage types in a comprehensible and accessible manner. To this end,

135      we expand previous works on data leakage by providing a comprehensive overview and easily

136      accessible visual representation of various leakage scenarios, categorized in a user-centric

137      and intuitive fashion. We hope that this overarching survey on data leakage scenarios

138      encourages more careful design and evaluation of ML pipelines and inspires further

139      investigation in this area. We aim to equip readers with the necessary tools to effectively

140      recognize leakage in their own (and others') work. This understanding will aid in avoiding these

141      pitfalls, fostering more robust and reliable ML-based analyses.

142      We would like to note that this work does not aim to cover the entire field of machine learning,

143      as it is too vast. For instance, we only touch upon time series analysis and do not address

144      unsupervised learning, which come with their own unique set of limitations and considerations.

145      We focus on supervised learning, however, most of the concepts and guidelines presented

146      here are generally applicable. The authors have noted the misconceptions and malpractices

147      discussed here in open-source code available on the Internet, as well as in code written by

148      themselves, students, or collaborators. These observations span various skill levels, ranging

149      from beginners to domain experts and data analysis experts. Therefore, the insights shared

150      here can provide guidance for everyone from novice to advanced ML practitioners,

151      researchers, reviewers, and editors.

152      We start with a brief introduction of ML basics and the cross validation (CV) procedure (section

153      2) that will serve as a guide to understand the concepts used in the rest of the article. This

154      section is divided in three parts: 2.a) ML concepts, 2.b) Cross-validation basics, and 2.c) Steps

155      while designing a ML pipeline. Next, we present various examples of leakage in ML pipelines

156      together with empirical examples and illustrations (section 3). Finally, we discuss possible

157      mitigations strategies (section 4) followed by general conclusions and key takeaways (section

158      5).

159

160

**2. Supervised Machine Learning: Pipelines and Evaluation**

2a. Supervised Machine Learning concepts

In a supervised machine learning task the user has access to labeled data consisting of $n$ feature-target pairs $S = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), \ldots (x_n, y_n)\}$ where $x_i$ are the features and $y_i$ are associated targets. The data samples are assumed to be independently and identically distributed (I.I.D.) and sampled from a fixed probability distribution. The task of a ML algorithm is to learn a function or a model that maps features to a target; $f(x_i) = y_i$. A model with discrete output is called a classifier, while one with continuous output is a regressor, two commonly encountered scenarios. The goal is to learn a model that generalizes on unseen data by providing accurate predictions. A model is composed of parameters (e.g., weights in a multiple linear regression) and often includes hyperparameters (e.g., regularization parameter $\lambda$ of ridge regression). Both contribute significantly to a model's ability to generalize. While the parameters are learned from the data using an optimization procedure, typically involving empirical risk minimization, the hyperparameters need to be either set by the user or "tuned" by searching for values that yield accurate predictions on hold-out data.

*2.b Cross validation basics: model assessment and model selection*

The goal of ML is to create models that accurately predict outcomes on unseen data, which requires learning generalizable information. However, because real-world test data (e.g., future patients or scenarios not yet encountered by a self-driving car) are typically not available, ML practitioners often hold out a portion of the available data as a proxy for test data to evaluate a model's generalization performance. Assuming that the underlying probability distribution of the data does not change, such an estimate helps with *model assessment* as an indicator of what to expect on new data.

Cross validation (CV) is frequently employed for *model assessment* (Fig. 1) as it makes efficient use of available data (33–36). In a $k$-fold CV scheme, the data is divided into $k$ non-

186 overlapping equally sized sets or folds. In each iteration of the CV procedure, one of the folds

187 is used as the test data, while the rest are used for training. Iterating through all folds

188 completes one CV run (also called a repeat). The average performance across all folds is

189 computed as an estimate of generalization performance. To minimize biases that could arise

190 due to data splitting, it is a standard practice to repeat the CV process multiple times with

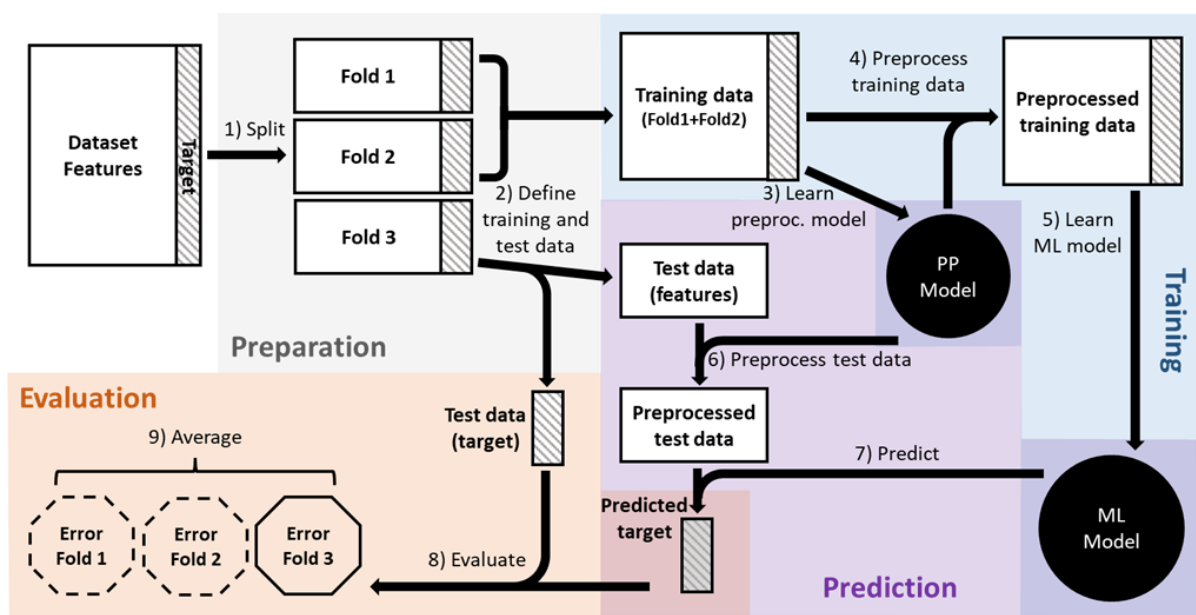191 different splits (e.g. 5 times repeated 5-fold CV) (37).

192



194 **Figure 1**: **A schematic representation of the cross validation (CV) scheme:** Here, we

195 illustrate a single repeat of a $k$-fold CV with three folds ($k = 3$) with the third fold being used

196 as the test data.

197

198 CV is also employed for *model selection* to select a model from a set of competing options

199 One example of these competing options are different models arising from hyperparameter

200 tuning (38) such as the cost of an SVM (support vector machine) (39). Another example are

201 different models arising from pipelines employing different preprocessing and/or learning

202 algorithms. The model with highest generalization performance is typically selected. CV

203      provides a general and practical method for model selection even in the case of complex

204      model parametrization. This is often found in ML algorithms where model selection statistics

205      like Akaike's information criterion might not be feasible. For a more comprehensive coverage

206      of the topic, we refer the reader to excellent sources (e.g. (38,40,41)).

207      Both *model assessment* and *model selection* are often a part of a ML-based data analysis

208      pipeline. However, as will be discussed in more detail below (see Section 3.a), problems arise

209      if the two roles are confused (37). Therefore, to cleanly and explicitly differentiate between

210      these two roles of CV (model selection and estimating generalization error), it is necessary to

211      use nested cross-validation (42), also known as double cross-validation. Within a nested CV

212      scheme, the inner CV encompasses all data-dependent decisions and performs model

213      selection (e.g., determining optimal hyperparameters or feature selection) while the outer CV

214      is responsible for model assessment, (i.e. evaluating the model after a finalized model

215      selection on previously completely unseen new data). The key point here is that any decision

216      made on data (i.e. the decision to select a specific model) requires yet again more data, that

217      was not involved in making the decision, to correctly estimate the generalization error.

218      *2.c Designing a ML pipeline*

219      The process of designing a ML-based data analysis pipeline can be broadly categorized into

220      the following steps: S-I) Task definition, S-II) Data collection and preparation, S-III) Data

221      preprocessing, S-IV) ML algorithm definition, and S-V) Definition of evaluation scheme and

222      metrics. If the goal of the analysis extends beyond assessing generalization performance,

223      additional steps might be employed, S-VI) Interpretation and deployment. While each of these

224      steps requires multiple decisions that must be made in a data-driven fashion, it is possible and

225      indeed necessary to define how each decision should be made a priori. Mistakes in the data-

226      driven decision making process can lead to data leakage. For more elaborate analysis

227      scenarios, we refer the reader to the CRoss-Industry Standard Process for Data Mining

228      (43,44).

229  *S-I) Task definition*: Definition of the target variable $y$ (e.g., disease status or behavioral

230  scores) and the features to be used (i.e., $x$, e.g., pixel values in images or functional

231  connectivity derived from neuroimaging data). Consideration of any confounds that can

232  obscure the intended feature-target relationship must be taken into account (e.g., age or sex

233  are often considered as confounds in biological and clinical applications).

234  *S-II) Data collection and data preparation strategies:* Here decisions need to be made both

235  before and after data collection. Before data collection, decisions to deal with known biases

236  should be made (e.g. equal sampling of males and females, or of case and control

237  observations; see (45,46) for a detailed treatment of this topic). After collection the data might

238  need preparation. This may involve subsampling (such as selection of only females for sex-

239  specific analysis), feature extraction like connectivity from brain imaging data, and feature

240  preparation like normalization of images. Importantly, we define data preparation as

241  processing exclusively applied to a single data point or sample independently of others. Open

242  and already prepared data are often available and are used directly by many practitioners.

243  S-III) *Data preprocessing strategy*: Optional data preprocessing steps involving

244  transformations applied across multiple samples are defined. These steps are typically applied

245  to the features, and may include feature normalization, feature selection, dimensionality

246  reduction, and treatment of missing values. Note that domain-specific data preparation and

247  (pre)processing is often employed and the reader is requested to refer to appropriate literature

248  for details.

249  *S-IV) ML algorithm definition:* One or more ML algorithms suitable for the task at hand must

250  be selected, such as classification for predicting disease status, or regression for predicting

251  continuous behavioral scores. That is, practitioners should a priori define a set of candidate

252  models to involve in model selection. If the ML algorithm includes hyperparameters, the

253  practitioner must either set the hyperparameter values or define a search space and search

254  strategy for tuning them using data (in the model selection process).

255 *S-V) Definition of evaluation scheme and metrics:* An evaluation scheme must be chosen for

256 model assessment, such as train-test split, $k$-fold CV or use of data to be collected in the

257 future. If the pipeline requires hyperparameter tuning, the chosen scheme should take this into

258 account, for instance by using nested CV. Evaluation metrics appropriate for the task must be

259 selected, such as classification accuracy and area under the receiver operating characteristic

260 curve (AUROC) for classification or mean absolute error (MAE) and coefficient of

261 determination ($r^2$) for regression.

262 S-VI) *Interpretation and deployment:* The selected model can be used to gain insights into the

263 structure of the data. Implicitly interpretable models provide parameters that can be used, e.g.,

264 weights of a linear SVM, or additional processing during or post model construction might be

265 needed, e.g., feature importance scores. In real-world application scenarios, the selected

266 pipeline is deployed for making predictions on new samples. In this case, the practitioner must

267 define how the new samples will be acquired and processed before making predictions. While

268 deployment is not considered in typical research settings, as we shall see, it serves as a useful

269 concept for avoiding some potential pitfalls.

270 **3. Leakage in ML pipelines**

271 Data leakage is a common and critical error in ML pipelines. Any data-driven choice made

272 within any step of a ML pipeline (see Section 2c), whether concerning preprocessing, learning,

273 or prediction, must be validated using new unseen data. Failure to use unseen data amounts

274 to leakage and results in inaccurate generalization performance estimates on the data at hand.

275 We take a general view of leakage to cover inappropriate use of data in different parts of a ML

276 pipeline which can lead to erroneous (either optimistic or pessimistic) estimation of

277 generalization performance or results in non-deployable models. Below, we describe several

278 types of leakage assuming that the ML pipeline employs CV for estimating generalization

279 performance.

280

281 *3.a Test-to-train leakage*

282 We begin with a type of leakage which we call test-to-train leakage as in this case information

283 is leaked from the test set into the training process. Several scenarios can lead to test-to-train

284 leakage, such as failure to separate training and test data, failure to consider dependent

285 samples, application of preprocessing before data splitting, and improper model selection.

286 The most straightforward case happens when the separation between training and test data

287 is not followed (47), i.e. the test samples are used for training (Figure 2). As the model can

288 learn patterns in the test data, it can result in high test accuracy. However, this cannot be

289 considered a correct estimate of generalization performance as the test data was not unseen.

290 That is, when the separation between training and test data is breached, the model risks fitting

291 to the specific patterns present in the test set. For instance, a k-nearest neighbors (KNN)

292 model (with k=1) will simply remember all the training samples it has seen, and therefore will

293 achieve perfect prediction if the model is "tested" on previously seen training samples. Of

294 course, the resulting error estimate can't possibly hold on truly unseen new data, so that this

295 error estimate is overly optimistic. Another example is optimizing a model to predict a particular

296 test set. This has happened previously in a ML competition where multiple evaluations on the

297 test set were performed, leading to disqualification of the team and consequently withdrawal
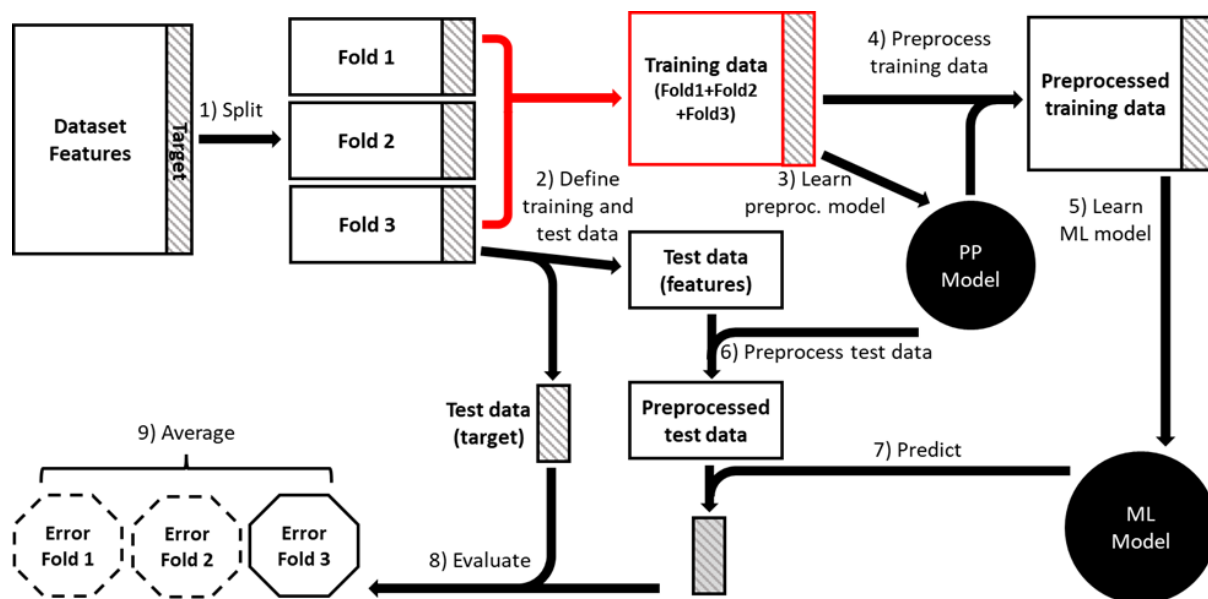
298 of the associated paper[1].

299

---

[1] https://dswalter.github.io/machine-learnings-first-cheating-scandal.html

13

300

**Figure 2**: **Test-to-train leakage:** Using the test data as a part of the training data leads to leakage, since the model is able to learn patterns from the test data during training, which usually results in overoptimistic generalization performance estimates. Red color indicates the problematic steps.

Another case where this type of leakage can happen is when the samples are not independent of each other. When data samples can be assumed I.I.D., randomly splitting them into training and test sets is sufficient. However, when the I.I.D. assumption is violated, more care must be taken to avoid related samples being split across training and test sets. An example of violation of the I.I.D. assumption is data that contains twins or siblings as is the case with the Human Connectome Project - Young Adult cohort (48). Functional Connectivity (FC) as a marker of brain organization is often used as a feature set to predict target variables such as behavioral scores in brain imaging research (49,50). Due to heritability, similarity between FC features (and likely also of target variables) is typically higher for twins and siblings than for independently drawn samples (51). If the twins are allowed to split between training and test sets, this is akin to duplicating the samples albeit noisily, therefore a model can learn about samples in the test set from their siblings in the training set. The prediction performance is

318    therefore higher when splitting family members across folds than if the family members are

319    grouped together (Figure 3) (29). This is an important pitfall, since most often the goal is to

320    build models that will generalize beyond specific families. The ungrouped CV cannot give us

321    an estimate of the error that we would expect for completely new samples, i.e. from new,

322    unseen families. Further examples demonstrating such leakage include, existence of the

323    same genomic loci in the training and test sets when performing cross-cell type predictions

324    (52), and using 2D slices from the 3D brain images of the same individual for training and

325    testing for predicting neurodegenerative disease (53).

326    In some task domains data does not follow the I.I.D. assumption (54). This happens for

327    instance in time series forecasting where the goal is to predict the future using historical

328    data (54,55). This happens because the data at consecutive timepoints are associated.

329    Hence, application of standard cross-validation is inadequate here as it disrupts the

330    temporal sequence by splitting sections of the time series and randomly assigning them

331    to training or test folds. Subsequently, past and future data is used inconsistently, i.e.

332    future data is used to predict the past, leading to leakage. This can be seen as a form of

333    test-to-train leakage, producing misleading estimates of predictive performance which will

334    not be representative if such a model is deployed in the real-world. To address this,

335    specific techniques, such as use of out-of-sample (holdout) test data corresponding to the

336    future (with respect to the training set), must be employed to obtain proper generalization

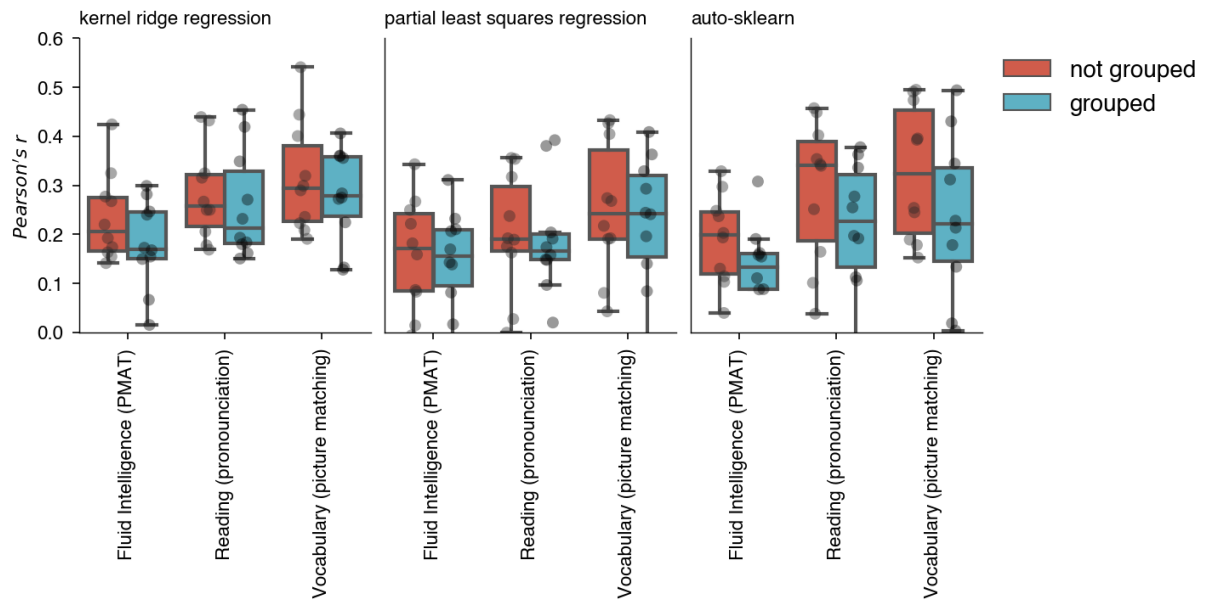337    estimates (54). We refer the reader to (16) for further details.

**Figure 3: Test-to-train leakage due to violation of I.I.D. assumption in cross-validation**:
Functional connectivity was estimated using resting-state functional magnetic resonance imaging data of the Human Connectome Project - Young Adult cohort (HCP-YA). Functional connectivity was then used to predict three psychometric targets (x-axis), each in a 10-fold cross-validation scheme. HCP-YA data contain siblings, and siblings are known to have similar connectomes. Therefore, allowing the siblings to be split across training and test sets (red bars, without grouping) leads to leakage while grouping siblings in training or test sets (blue bars, with grouping) shows overall lower accuracy (Pearson's r between true and predicted target, y-axis).

In addition to leakage due to the same or similar samples, data leakage can also happen via modeling preprocessing parameters (preprocessing leakage). A common case of such test-to-train leakage arises when data preprocessing, such as dimensionality reduction (e.g. principal component analysis (PCA)), confound removal, feature normalization or scaling, and imputation for filling in missing values, is applied to the whole dataset before splitting it for CV (16,29). Practitioners may not immediately recognize this as leakage, since the ML model is trained after splitting the data. However, estimating the preprocessing parameters on the

356     whole dataset invalidates the train-test separation (Figure 4). That is, data in the training set

357     is transformed dependent on data in the test set, and crucially, ML models can exploit this to

358     learn about the test set. Therefore, the resulting estimate of generalization performance is

359     likely to be overly optimistic, though it can also decrease the performance of the models (29).

360     Empirical demonstrations of such leakage in the literature include performing confound

361     removal (29) or feature selection on the whole dataset (Figure 5) (23,29,56), and oversampling

362     to counter data imbalance (25). It should be noted that such leakage can happen when

363     preprocessing either the features or the target values. For instance, when the target is created

364     by combining multiple variables (e.g., several behavioral measures) using a dimensionality

365     reduction method such as PCA.

366     This case of test-to-train leakage can be avoided by learning the preprocessing parameters

367     on the training set and then applying them to the training and the test sets. For example, if one

368     wants to apply feature selection, one should select features based on the training set after

369     splitting data. Crucially, this means that each iteration of CV will involve its own feature

370     selection process which may select a different set of features. Since this implies more

371     computation and also makes interpreting the results more difficult given that it is necessary to

372     keep track of different features in different CV iterations, practitioners sometimes erroneously

373     avoid splitting the data before data preprocessing. Note that data preparation strategies that

374     rely on a single sample and thus preserve the train-test separation do not lead to such leakage.

375     For instance, data imputation can be performed in a within-sample fashion such that missing

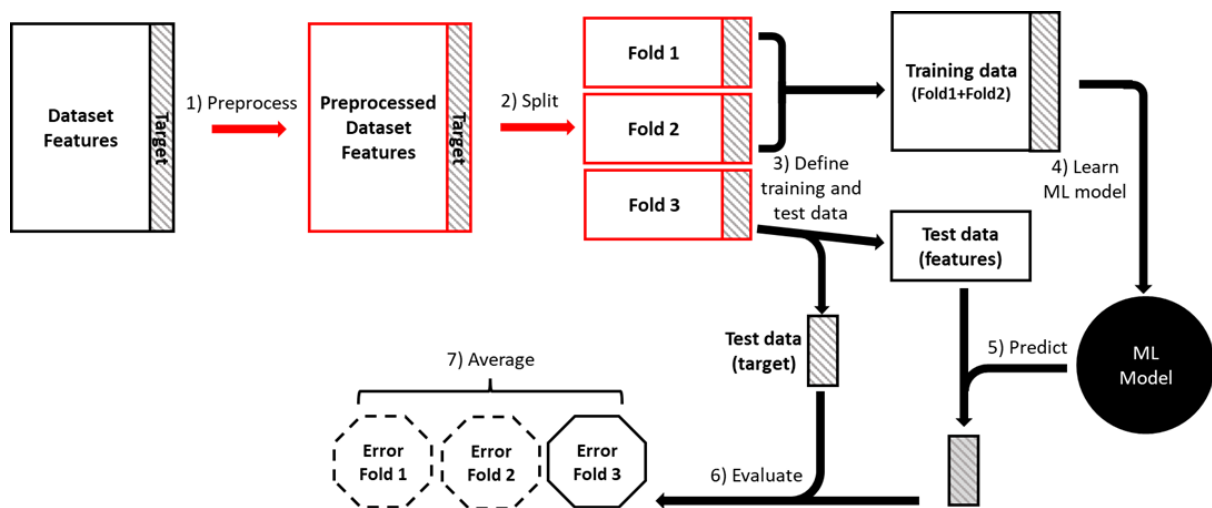376     values are estimated using other features of that sample (see, e.g., (57)).

377

378

**Figure 4**: **Test-to-train leakage:** Preprocessing is performed on the whole data before data

splitting, yielding a preprocessing model learned from both training and testing data. In a

correct implementation, the parameters of a preprocessing model should be estimated in the

training set and applied to both training and test sets. Red color indicates the problematic

steps.

```
# leakage
features_to_use  = feature_select(X, y)
CV_error = run_cross_validation(X[:, features_to_use], y, k=5)

# no leakage
for k in 1:5
    train_idx = training_set_indices(k)
    test_idx = training_set_indices(k)
    features_to_use = feature_select(X[train_idx, :], y[train_idx])
    model = train_model(X[train_idx, features_to_use], y[train_idx])
    CV_error[k] = test(model, X[test_idx, features_to_use])
```

**Figure 5: Pseudocode for leakage-inducing and leakage-free feature selection.**

A particular case of test-to-train leakage occurs when conflating the two roles of CV, model

assessment and model selection (see Section 2.b). Although not immediately obvious, this

can be considered test-to-train leakage, because by running CV for many different models,

391    and subsequently making a decision based on the results (i.e. selecting one model among

392    them), the test folds used in the CV essentially become part of the training data. It is important

393    to highlight here that training data does not just refer to data that an algorithm uses to fit

394    parameters but also to data that researchers use to make data-driven decisions. Therefore,

395    the error estimate from the model selection CV is not a valid estimate of true generalization

396    performance. For instance, consider an ML algorithm with a single hyperparameter such as a

397    linear kernel SVM with its hyperparameter called cost within a 5-fold CV. For each fold the

398    cost value that provides the lowest error on the set is used. The CV estimate of error is

399    calculated by averaging across the test sets. As a result, the error obtained during model

400    selection is likely an overoptimistic estimate of the generalization error (model assessment)

401    (23,24,37,58). Using nested CV in which the hyperparameter is tuned in an inner CV and the

402    selected model is applied on the test set avoids such leakage.

403    One of the questions a practitioner may face in this regard concerns the fact that different

404    models and hyperparameters are selected in each fold of the CV. Students wonder how they

405    can report and use that model, since there is no such thing as "the model". However, this point

406    of view fails to acknowledge that simply fitting parameters of a model in each iteration of CV

407    will also always result in different models. Instead, nested CV presents an opportunity rather

408    than a challenge, since researchers can then easily test the stability of fitted parameters and

409    hyperparameters chosen in the model selection process over multiple iterations by inspecting

410    each trained (and selected) model. We provide an empirical illustration, again using

411    neuroimaging data to predict behavior, that indeed shows that CV estimates are overoptimistic

412    compared to nested CV estimates (Figure 6).

413    Such leakage is not restricted to hyperparameter tuning and can happen with any data-driven

414    choices, such as selection of an algorithm (e.g., SVM versus random forests) and data

415    transformations (e.g., PCA, univariate feature selection) (37). Such choices should be treated

416    in the same way as hyperparameters, i.e. tuned and evaluated using nested CV. In other

417    words, all data-driven choices within a ML pipeline should be considered as a part of learning,

hence they must be validated on data not seen by the models to correctly estimate
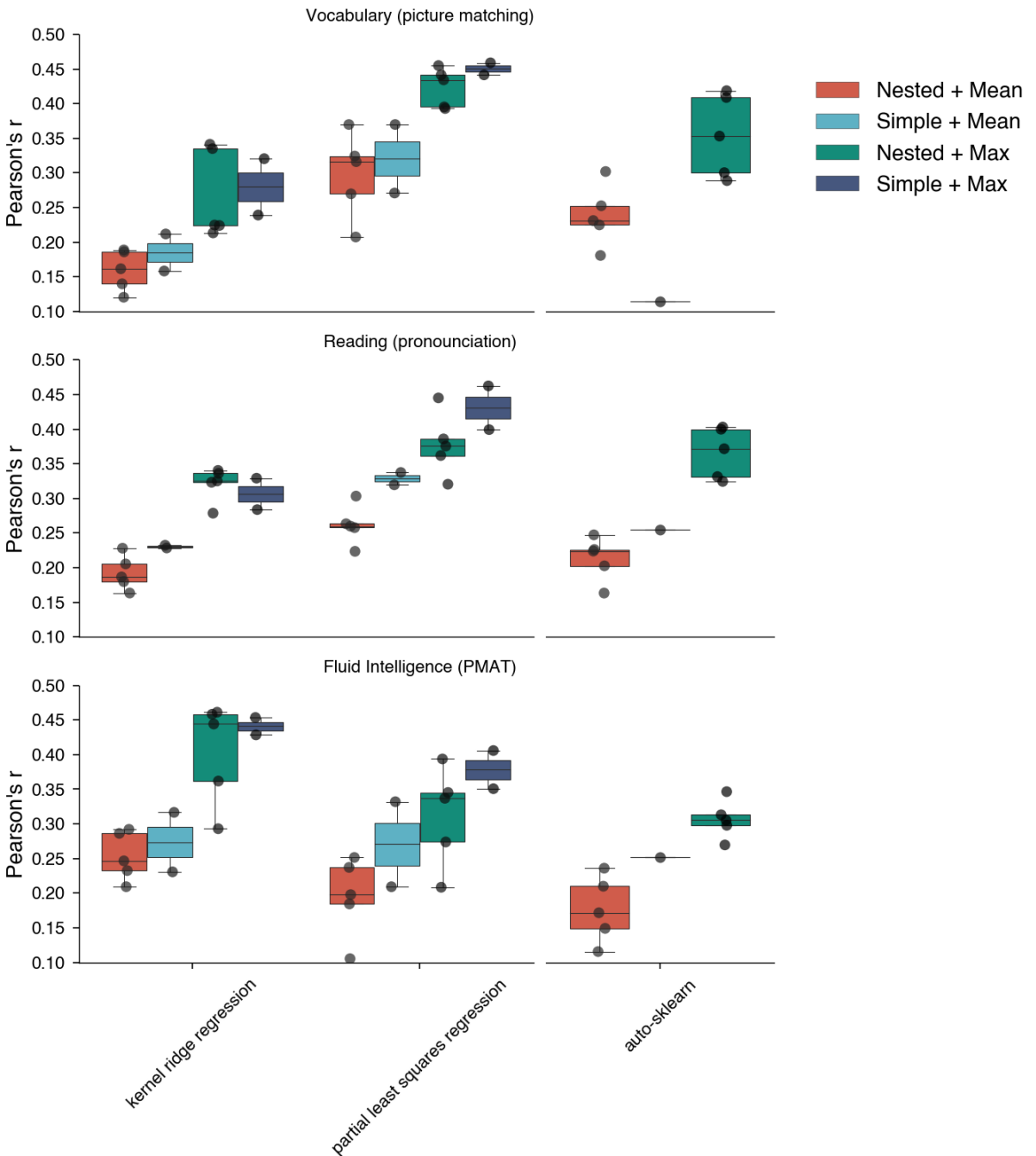
419 generalization performance. For more details see (59).

420



421

**Figure 6: Effect of erroneous usage of CV for estimating generalization performance:**

423 Three targets (behavioral scores) were predicted in a 5-fold CV with 5 repeats using a

424 subset of the HCP-YA S1200 release consisting of 369 unrelated subjects (192 males, 177

425    females) with ages ranging from 22 to 37 (*M*=28.63, SD=3.83). In one case the CV was used

426    for hyperparameter tuning and estimation of generalization performance simultaneously (i.e.,

427    "simple" CV) and the mean values (light blue) and maximum values (dark blue) across folds

428    were reported. In the other case nested CV was performed such that hyperparameter tuning

429    was performed for each test fold independently in an inner CV loop applied on the training

430    folds of the outer CV loop. Again, mean values (red) and maximum values (green) across

431    folds are reported. The effect of reporting the maximum value across folds rather than the

432                                    mean can be visualized.

433

434    *3.b Test-to-test leakage*

435    Test-to-test leakage is a covert type of leakage that arises due to erroneous information

436    sharing between the test samples. For the majority of ML applications samples in the test set

437    should be treated independently. That is, processing and prediction of a given test sample

438    should not depend on information from other test samples. Test-to-test leakage happens when

439    the test set is used for estimating preprocessing parameters (Figure 7). In other words, instead

440    of a single preprocessing model derived from the training samples (see Figure 1 for correct

441    implementation) two preprocessing models are estimated: one using the training samples and

442    applied to training samples, and another using the test samples and applied to test samples.

443    For example, a practitioner may wish to demean their features. They may then estimate the

444    mean of each feature in the training data and use these estimations to demean the training

445    set. The error occurs if instead of also applying these estimations to the test data, the

446    practitioner then estimates the mean values on the test data to demean the test data. This is

447    wrong, because it implies that individual samples in the test data are demeaned depending

448    on other samples in the test data.

449    Another complication that arises in a pipeline that does not treat an individual test sample

450    independently of other test samples, is that it cannot be deployed, and it might fail completely

451    when applied to a single test sample. In the previous example of demeaning for instance, the

452    demeaning could not be applied since the demeaned feature will be zero. In addition, these

453    kinds of preprocessing steps attempt to estimate the parameters for a population, and

454    estimates of parameters such as the mean are unlikely to be accurate from the typically

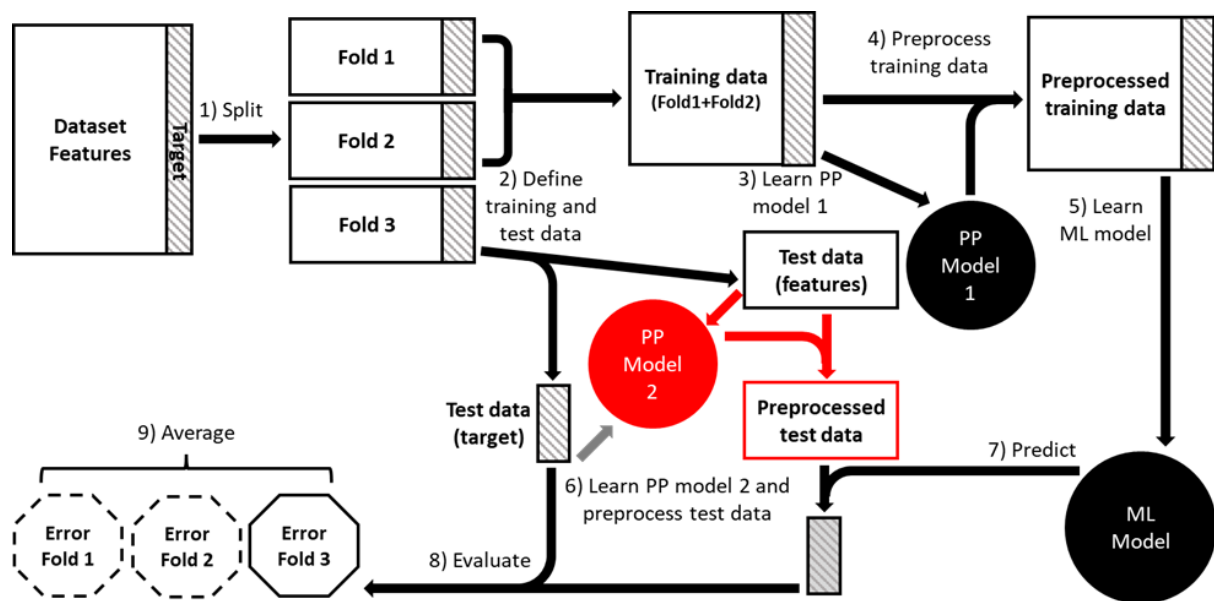455    smaller test set even if they did not result in impermissible operations.

456

457    

458    **Figure 7: Test-to-test leakage:** Preprocessing is performed in the test and training sets

459                        independently. Red color indicates the problematic steps.

460

461    Another case of test-to-test leakage can happen due to inappropriate use of the predictions

462    obtained in CV. With repeated CV, which is usually recommended to avoid biases due to

463    random data splitting, one obtains multiple predictions for each sample. Combining each

464    sample's predictions across CV repeats, e.g., by averaging, causes leakage (Figure 8).

465    Although this may seem like a rather elegant way to obtain a single prediction per sample and

466    in turn a single error estimate, this procedure has two negative consequences. First, it

467    generates ensemble results, and the performance reflects an ensemble of models built across

468    repeats and not a single model as the practitioner intended and might claim. Second, it

469    increases the effective number of CV folds ($k$) because across CV repeats the same data

470    point is predicted using different combinations of training samples, effectively reducing the

471    out-of-sample data points. With a high number of repetitions, the averaged prediction would

472    be influenced by all other data points akin to leave-one-out (LOO) CV, and the claim that the

473    results reflect $k$-fold strategy would be wrong. These two reasons can lead to overoptimistic

474    results.

475    We note that LOO per se is not problematic if properly implemented and with a correct

476    evaluation metric that can be calculated for each sample separately and does not combine

477    the samples. For instance, classification accuracy is suitable in a classification task, whereas

478    AUROC is not. Similarly, when using LOO mean absolute error is appropriate but Pearson

479    correlation should not be used.
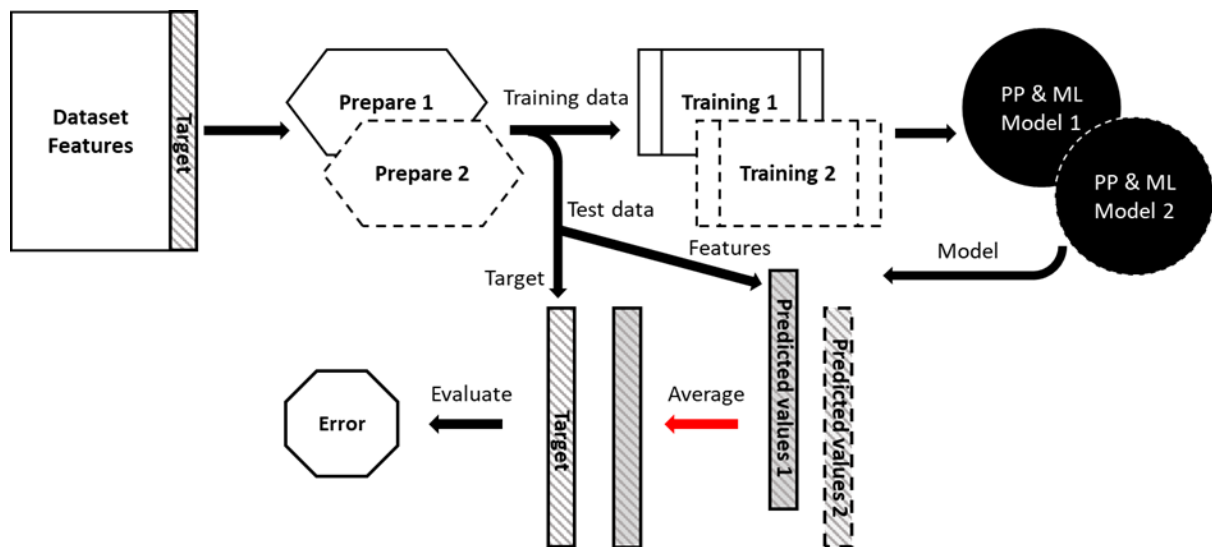
480



481

482

483    **Figure 8: Test-to-test leakage due to averaging of predictions:** For simplicity we show

484    two runs of CV, and the major parts of CV are abstracted. Red color indicates the

485    problematic steps.

*3.c Feature-to-target leakage*

487    Feature-to-target leakage occurs when the target is constructed using one or more features

488    in the first place (Figure 9). For example, data is first clustered using some or all the features

489    and the resulting cluster IDs are used as the target. Subsequent generalization estimation

490    using CV on this data is likely optimistic as the supervised classification algorithm will simply

491    need to reverse engineer the clustering process which it should be able to in most cases. Note

492    that it is valid to use cluster IDs as target to train a classification model and apply it to new
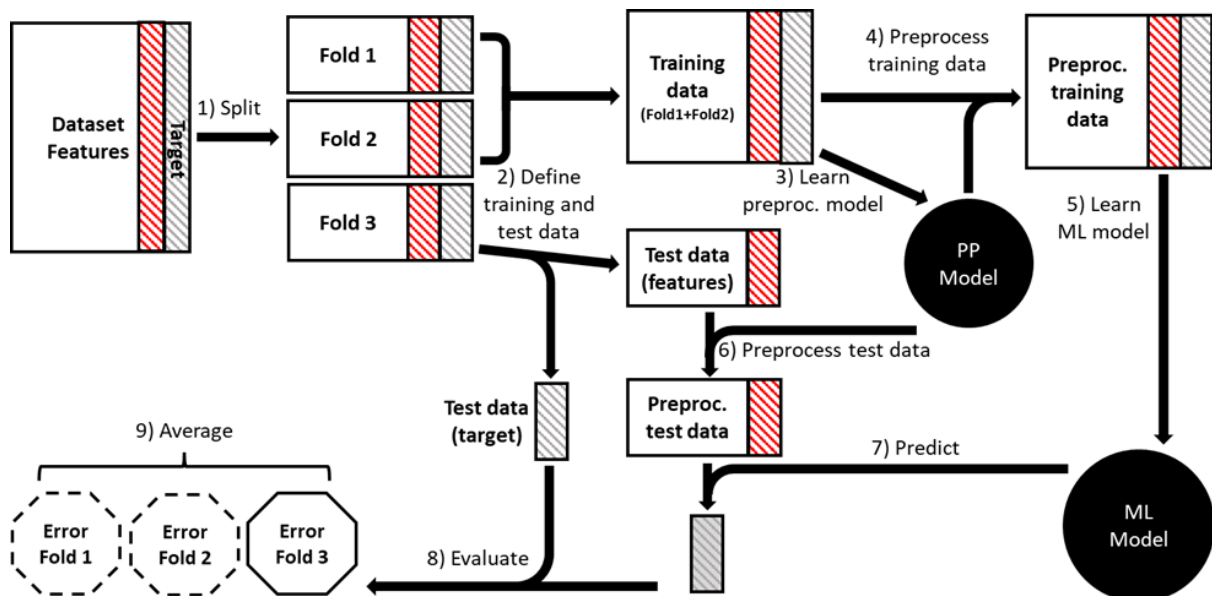
493    unseen data.

494

495



496    **Figure 9**: **Feature-to-target leakage:** The features include a variable, indicated with red

497                color, that directly contributes to the target.

498

499    Feature-to-target leakage is also evident when one or more variables informative of the target

500    are available during learning but not at the deployment phase. Since these features carry

501    information that the models can learn from, it will lead to a high CV accuracy and the model

24

502  might seem acceptable. However, this model is not deployable as the crucial features will not

503  be available at the test time. As an example, consider a healthcare application where the goal

504  is to diagnose a disease based on radiology data. To generate training data, radiology experts

505  manually label each image as diseased or healthy. While labeling the data, the experts also

506  make notes such as the size and location of the abnormality. If these notes are used together

507  with the radiology images when building and validating models, it can lead to above mentioned

508  issues. Given that these notes are predictive, the models will utilize them, and as a result, the

509  CV will likely demonstrate high performance. However, these manually created notes will not

510  be available during real-world use as the goal of such an application is to avoid manual

511  annotation. Therefore, a model trained in this way cannot be deployed.

512  More generally, this type of leakage occurs due to differing train and real test data distributions

513  which again violates the I.I.D. assumption. However, as CV is performed only in the training

514  data the predictive information available is used providing an optimistic estimate of

515  generalization performance. As such, this scenario is similar to reliance on confounds (60)

516  and shortcut learning (61) where a model learns unintended signals prominently present in

517  the training data which then hampers generalization to unseen data in which these signals are

518  missing. While this scenario could be viewed as learning incorrect information rather than

519  leakage, we categorize it as leakage because it results in overestimation of the model's ability

520  to generalize which is typically the culmination of academic exercises. Note that feature-to-

521  target leakage has also been termed as target leakage[2] but we reserve this term for another

522  type of leakage described below.

523

524

---

[2]http://downloads.alteryx.com/betawh_xnext/MachineLearning/MLTargetLeakage.htm
https://h2o.ai/wiki/target-leakage/

525   *3.d Target leakage*

526   We define target leakage as the scenario where the target is utilized in the prediction process.

527   This is an overt kind of leakage that can occur due to use of specific algorithms and

528   preprocessing steps that rely on the target.

529   An example is misapplication of the partial least squares (PLS) algorithm, which is popular in

530   several fields. PLS is a bilinear factor model that identifies a new space representing the

531   covariance between the features $X$ and target $Y$ spaces. As PLS estimates a shared latent

532   space, it generates beta values for both $X$ and $Y$. If the beta values $Y$ are also used instead of

533   only that of $X$ to make predictions, then one requires the target values of the test examples,

534   resulting in data leakage. When performing CV, the target values of the test data, are available

535   and thus one could provide them at the test time to generate the latent space of the test data,

536   but this will result in leakage in addition to creating a non-deployable model. Thus, when using

537   PLS the practitioner should make sure that only $X$ values are used for prediction. Note that

538   most libraries provide a correct implementation of PLS.

539   Another case of target-leakage can occur while trying to mitigate measurement bias by means

540   of data harmonization. As data collection becomes more accessible and widespread, pooling

541   data together from different sites has become increasingly common. Differences in data

542   collection protocols and measurement devices can induce systematic biases. To address this,

543   batch-effect removal (22) or data harmonization (62) is used. If the data are harmonized

544   across sites (e.g., by matching covariances) too aggressively, it could also remove variance

545   of interest that is related to the target. To avoid this, harmonization methods can preserve

546   variance related to user-specified covariates which can also include the target. When

547   harmonizing new data, the same covariates must be provided for each test sample. Effectively,

548   while it is beneficial to preserve variance related to the target while harmonizing, the target

549   values must be available when making predictions on test data. This presents two possible

550   ways to implement harmonization while estimating generalization performance using CV, both

551    of which lead to leakage. First, one can harmonize all the data before creating splits which

552    violates the train-test data separation requirement. Second, the target values of the test set

553    are needed for harmonization which amounts to target leakage and also precludes

554    deployment. Overall, the use of harmonization in a ML pipeline while explicitly preserving the

555    target-related variance can be considered as target leakage.

556    *3.e Dataset leakage: Dataset decay*

557    In traditional statistics it is widely recognized that testing multiple hypotheses (e.g., mass-

558    univariate hypothesis testing) can result in spurious discoveries (type I error), and hence

559    approaches to correct for multiple comparisons are utilized (40,63,64). The multiple testing

560    issue might not seem pertinent to an individual researcher who only aims to test a single

561    hypothesis using a given dataset. On a larger scale, however, multiple researchers testing

562    different hypotheses on the same dataset leads to an increased global likelihood of false

563    positive findings, referred to as "dataset decay" (65). In other words, dataset decay means

564    that as more hypotheses are tested in a specific dataset, the utility of such a dataset to yield

565    generalizable results decreases.

566    In some fields, for example neuroimaging, the availability of large datasets suitable for ML

567    analysis is relatively limited, e.g. Human Connectome Project and Alzheimer's Disease

568    Neuroimaging Initiative. Consequently, these datasets are extensively used by a multitude of

569    researchers in the field. This widespread usage has sparked a competitive environment where

570    there is a continual race to develop new methods that outperform existing state-of-the-art

571    scores. However, this intense focus on a few datasets and high accuracy can lead to dataset

572    decay such that the overuse and repeated analysis of the same datasets increases the rate

573    of false positive findings. This, in turn, would diminish the effectiveness of such models for

574    future research and innovation, decreasing generalizability of the findings. For this reason, we

575    consider this issue as a type of leakage (dataset leakage).

576    Furthermore, there is another angle to dataset leakage, which happens when a researcher's

577    data analysis strategy is informed by previous analyses carried out on the same data, i.e.

578    adaptive data analysis or circularity (66). This not only increases the number of comparisons

579    but also introduces dependency between analyses, frequently leading to further false

580    discoveries and inflated performance estimates (67). In fact, our empirical demonstrations

581    may be a good example of such overfitting: In both Figure 4 and Figure 6, kernel ridge

582    regression and partial least squares regression were used for predicting a number of

583    behavioral scores based on functional connectivity. These two models are well known to work

584    on this task and have been chosen for this exact reason (50,68). Importantly, they are known

585    to work well on the exact dataset we have used here - the HCP-YA. This could explain why

586    instances of these models outperformed models found by Auto-ML. However, it is important

587    to note that this conclusion for this specific case is only suggestive and further evaluations

588    using additional datasets are needed to assess whether the algorithms and resulting outcomes

589    are truly overfitted.

590    <u>3.f Confound leakage</u>

591    A ML model can be influenced by confounding factors which in turn can influence its

592    predictions. Confounding factors are related to both the features and the target. Depending

593    on the goal of a study, it may or may not be desirable to minimize their impact. If the goal of a

594    study is to simply predict a target as well as possible, then confounding may not necessarily

595    be a worry (69). But if researchers want to gain insight about the specific relationship between

596    features and target independent of the confounding factors, strategies to mitigate the effect of

597    confounders must be applied. A common example of a confounding factor in brain imaging

598    studies is age. For instance, brain imaging can accurately reflect a person's age and can also

599    contain information about age-related diseases (70). The problem here is that the model

600    simply learns how brain images may change with the natural aging process but may not learn

601    anything about the specific changes and processes related to pathology. From a prediction

602    point of view this may be an issue in some cases, for instance a young person with pathology

603 will not be identified as a patient. Furthermore, such a model is likely less helpful in gaining

604 specific biological insight than a model that learns about processes specifically involved in the

605 disease.

606 Thus, if not properly controlled or accounted for, confounding can lead to correct predictions

607 but may mislead the researcher to incorrectly conclude that their model is learning about the

608 specific disease-related processes. For example, a particular feature of brain structure might

609 be erroneously deemed important in the prediction of Parkinson's disease, whereas the

610 association is actually due to the confounding effect of age. In other words, the brain structural

611 feature changes with increasing age, and increasing age also leads to a higher likelihood of

612 developing Parkinson's disease, but the structural feature in reality has nothing to do with the

613 disease.

614 However, even if ML practitioners do not care about interpretability or insight with respect to

615 their model and how it represents the relationship between features and target, confounding

616 factors can be a problem. That is, confounding can also degrade the model's performance,

617 especially when the model is deployed in environments where the distribution of the

618 confounding variable is different. This is akin to the assumption that the data used to train a

619 model should follow the same distribution as the data for which the model is deployed, with

620 the only difference being that confounds in this scenario are not modeled explicitly (but

621 importantly change the relationship between features and target). Thus, the model trained and

622 evaluated in a particular distribution (of the confounding factor) will lead to overly optimistic

623 estimates of the generalization error, that may not tell us about how the model will perform on

624 data for which the distribution of the confounding factor is changed.

625 Feature-wise confound removal by means of linear regression (i.e. confound regression) is a

626 standard method used to deal with confounding effects in retrospective data analyses which

627 is a common scenario in ML (71,72). As mentioned before, it is recommended to perform

628 confound regression in a CV-consistent manner to avoid test-to-train data leakage (73).

629     However, it is important to highlight that the process of confound regression itself can leak

630     information into the features especially when the feature-target distribution is skewed (74). In

631     this case, variance associated with the confounds is injected into the features rather than being

632     removed as expected. This type of leakage becomes particularly problematic when the

633     confounds are strongly associated with the target and the leakage increases with the number

634     of features. Confound-leakage can lead to above-chance generalization performance

635     estimates, even when the relationship between the features and the target is destroyed. This

636     provides a method to check for potential confound leakage by shuffling each feature

637     independently and performing CV. If the CV accuracy is above-chance, then one can conclude

638     that confound leakage is possible on this dataset.

639     **4. Possible mitigation strategies**

640     In this section we provide advice to improve common reporting practices to facilitate the

641     detection of leakage as well as to increase reproducibility in ML pipelines. First we should

642     mention that excellent recommendations exist for reporting ML models such as Model Cards

643     (75) and Data, Optimization, Model and evaluation (DOME) recommendations (76). Further

644     guidelines have been proposed for specialized domains such as biomedical applications (77)

645     and minimum information about clinical artificial intelligence modeling (78). Data processing

646     and quality reporting have also been discussed (e.g. Datasheets for Datasets (79) and Data

647     Nutrition Labels (80)).

648     However, as the application contexts and modeling intricacies of ML-based analyses expand

649     and grow in complexity, we think it is necessary to refine current recommendations, especially

650     by making data processing and model selection and assessment strategies more transparent.

651     Effective communication of ML pipelines can take various forms. Textual descriptions offer

652     high-level overviews of a pipeline and its components, providing a conceptual understanding.

653     However, such descriptions can fall short.

654    When documenting the experimental setup, it is crucial to provide comprehensive details.

655    However, ambiguous descriptions hinder understanding and replicability. For instance, stating

656    that "default (hyper)parameters" were utilized for a given algorithm is inadequate as default

657    hyperparameters are not universally defined and can vary between software implementations

658    and labs. Moreover, defaults may change over time and with software updates. Therefore, to

659    guarantee replicability, it is essential to report the specifics of the set up including data

660    processing, model hyperparameters, training and evaluation procedures, as well as the

661    software packages used and their corresponding versions.

662    In addition, there are several instances where the method description, whether written or

663    verbal, is insufficient to detect leakage. In some cases, the description seems to be correct,

664    even when leakage is present. Therefore, researchers and reviewers should not only rely on

665    the written description, but instead insist on reviewing the code written to implement the

666    methods. To illustrate, let us look at this example method description: "We performed feature

667    selection using LASSO followed by an SVM classifier with a radial basis function kernel. Both

668    feature selection and hyperparameter tuning were performed on the training folds within a

669    cross-validation loop". Based on this, a reader would assume that the procedures were

670    correctly implemented. However, when we were unable to replicate the results on the same

671    data, we decided to examine the code more closely. The issue became immediately apparent:

672    test-to-train leakage during feature selection using LASSO. Although both steps, feature

673    selection and hyperparameter optimization, were performed within a CV loop as mentioned in

674    the text, these two steps were implemented in two separate CV loops (see Figure 5). The first

675    CV loop performed feature selection, correctly only on the training folds, and noted which

676    features were selected (nonzero LASSO weights). The features selected more than once

677    across CV folds were retained. Then in another CV loop, an SVM classifier was trained on the

678    training folds with hyperparameter tuning while using only the selected features and tested on

679    the test fold. It is clear that this procedure causes leakage even though the text description

680    can be interpreted as being correct. It is crucial to note that the feature selection method (as

681 long as it is data driven) or the fact that the same CV fold structure was used for the two loops,

682 does not prevent the leakage. The problem arises in the first CV loop where choice of which

683 features to use was made using all the data. Therefore, the textual description of the process

684 was insufficient to detect the issue of leakage. However, a quick examination of the underlying

685 code clarified the situation. This highlights the importance of transparency in machine learning

686 research and practice, and we strongly encourage researchers and practitioners to share their

687 code. In addition, to avoid such errors before publication it is crucial that labs implement their

688 own standard procedures to perform adequate code review. Since internal code review might

689 not catch all errors, a greater emphasis should be put by journals and reviewers on reviewing

690 code during the peer-review process.

691 Sharing source code allows for a detailed examination of the implementation, enhancing

692 transparency and enabling replication or modification. While sharing their code may indeed

693 expose it to critical review, it is essential for progress in this predominantly software-driven

694 discipline (81,82). There may be errors spotted and improvements suggested, but this iterative

695 process of refinement is an integral part of scientific advancement. In this regard, we echo the

696 call to make research code openly available. This practice would aptly put importance on

697 correctness of the pipelines and ensuing analyses. The authors admit that they themselves

698 have done this with less stringency in the past, but we aim to do better. We also recommend

699 that ML practitioners, especially in early career stages, request code reviews to identify and

700 fix any issues with their implementation. Importantly, however, the responsibility should not be

701 off-loaded to early career researchers and standard procedures (to ensure good quality code

702 and research) need to be implemented at an organizational level. For example, each lab can

703 perform code reviews using one or more skilled reviewers. Such a code review should focus

704 on the correctness of the code itself and not on its functionality or whether it produces the

705 desirable output such as high accuracy. However, we recognize that it can be challenging to

706 find skilled individuals that can perform a proper review.

707 Openly sharing data, trained models, and other research artifacts/tools fosters reproducibility,

708 and encourages collaboration, enabling the broader community to validate and build upon the

709 work. However, it must be noted that the shared data should not be affected by data leakage

710 (e.g., during preprocessing). Detecting such instances of leakage can be challenging, if not

711 impossible once the data is shared.

712 **5. Conclusions**

713 Data leakage presents a significant challenge in machine learning. Identifying and preventing

714 leakage is essential for ensuring reliable and robust models. To this end, we have provided

715 several examples and detailed explanations of data leakage instances, along with tips on how

716 to identify them. Below, we present a few crucial points that underlie most of the leakage cases

717 we have presented in this article.

718 ● Ensure strict training-test set separation.

719 ● Ensure that performance metrics are calculated on truly unseen data that has not been
720    used anywhere in the pipeline previously.

721 ● Model selection and model assessment should be done with a nested CV.

722 ● State the goal of your ML pipeline: Search for a feature-target relationship, assess
723    generalization performance or deployment? Clarifying the goals early on will help
724    practitioners to design and implement a correct and appropriate pipeline.

725 ● While academic applications of machine learning often do not involve deploying
726    models, considering whether a pipeline can actually be used to make predictions on
727    genuinely unseen data can aid in identifying potential instances of data leakage. Of
728    note, check if features are available after deployment. Can a model be applied to future
729    test data not currently available? Can it be applied to a single test example?

33

- Detailed description of methods as well as sharing your code publicly is an effective way to ensure transparency in your pipeline design. In addition, releasing your models enables users to test data gathered after the model's publication, which further boosts confidence in the model's ability to generalize well.

- When possible, opt for using well-established software packages and libraries instead of creating standard procedures from scratch. We certainly do not discourage code implementation for learning purposes, but we recognize that code testing can be time-consuming and challenging. Hence, using standardized code in production environments is typically a more effective choice.

- Ensuring the correctness of a ML pipeline should take precedence over its output or even its replicability. A flawed pipeline might yield accurate and replicable results, but that should not be mistaken as an indication of the validity of the models or the ensuing results.

Finally, in addition to leakage, several other pitfalls and issues exist and deserve attention, (real-world usefulness of benchmark data (83), dataset biases (46,84) and deployment challenges (85)). However, it is not possible to cover all those aspects in a single paper. We recommend that readers stay vigilant and pay attention to issues that might affect their specific analysis set up.

**6. Ethics approval and consent to participate**

The ethics protocols for analyses of these data were approved by the Heinrich Heine University Düsseldorf ethics committee (No. 4039).

**7. Availability of data and materials**

Access to data of the HCP can be requested on ConnectomeDB (https://db.humanconnectome.org/app/template/Login.vm).

## 8. Acknowledgements and funding

## 9. Author's contributions

L.S.: Designed the experiments, developed code, performed the analyses, contributed to discussion and interpretation of results, and writing of the manuscript

E.N-S: Contributed to discussion of results and writing of the manuscript

J.D.: Contributed to discussion and interpretation of results, and writing of the manuscript

S.B.E.: Contributed to discussion and interpretation of results, and writing of the manuscript

M.G.: Contributed to discussion and interpretation of results, and writing of the manuscript

S.H.: Contributed to discussion and interpretation of results, and writing of the manuscript

V.K.: Contributed to discussion and interpretation of results, and writing of the manuscript

A.K.: Contributed to discussion and interpretation of results, and writing of the manuscript

J.L.: Contributed to discussion and interpretation of results, and writing of the manuscript

B.C.L.: Contributed to discussion and interpretation of results, and writing of the manuscript

780 **<u>Bibliography</u>**

781  1.  Jiang T, Gradus JL, Rosellini AJ. Supervised machine learning: A brief primer. Behav
782      Ther. 2020 Sep;51(5):675–87.

783  2.  Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al.
784      Generative Adversarial Nets. Advances in Neural Information Processing Systems.
785      2014;

786  3.  Sutton RS, Barto AG. Reinforcement Learning, second edition: An Introduction
787      (Adaptive Computation and Machine Learning series). second edition. Cambridge,
788      Massachusetts: Bradford Books; 2018.

789  4.  Bhaskar H, Hoyle DC, Singh S. Machine learning in bioinformatics: a brief survey and
790      recommendations for practitioners. Comput Biol Med. 2006 Oct;36(10):1104–25.

791  5.  Sun AY, Scanlon BR. How can Big Data and machine learning benefit environment
792      and water management: a survey of methods, applications, and future directions.
793      Environmental Research Letters. 2019 Jul 3;14(7):073001.

794  6.  Swain S, Bhushan B, Dhiman G, Viriyasitavat W. Appositeness of optimized and
795      reliable machine learning for healthcare: A survey. Arch Comput Methods Eng. 2022
796      Mar 22;29(6):3981–4003.

797  7.  Varoquaux G, Cheplygina V. Machine learning for medical imaging: methodological
798      failures and recommendations for the future. npj Digital Med. 2022 Apr 12;5(1):48.

799  8.  Douglas MR. Machine learning as a tool in theoretical science. Nat Rev Phys. 2022
800      Mar;4(3):145–6.

801  9.  Larrañaga P, Calvo B, Santana R, Bielza C, Galdiano J, Inza I, et al. Machine learning
802      in bioinformatics. Brief Bioinformatics. 2006 Mar;7(1):86–112.

803  10. Wilkinson J, Arnold KF, Murray EJ, van Smeden M, Carr K, Sippy R, et al. Time to
804      reality check the promises of machine learning-powered precision medicine. Lancet

805    Digit Health. 2020 Dec;2(12):e677–80.

806    11.    Chen J, Patil KR, Yeo BTT, Eickhoff SB. Leveraging machine learning for gaining

807          neurobiological and nosological insights in psychiatric research. Biol Psychiatry. 2023

808          Jan 1;93(1):18–28.

809    12.    Qiu J, Wu Q, Ding G, Xu Y, Feng S. A survey of machine learning for big data

810          processing. EURASIP J Adv Signal Process. 2016 Dec;2016(1).

811    13.    Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-

812          learn: Machine Learning in Python. The Journal of Machine Learning Research.

813          2011;12:2825–30.

814    14.    Kuhn M, Wickham H. Tidymodels: a collection of packages for modeling and machine

815          learning using tidyverse principles [Internet]. 2020 [cited 2023 Aug 25]. Available from:

816          https://www.tidymodels.org

817    15.    Chen S, Sedghi Gamechi Z, Dubost F, van Tulder G, de Bruijne M. An end-to-end

818          approach to segmentation in medical images with CNN and posterior-CRF. Med

819          Image Anal. 2022 Feb;76:102311.

820    16.    Kapoor S, Narayanan A. Leakage and the reproducibility crisis in machine-learning-

821          based science. Patterns (N Y). 2023 Sep 8;4(9):100804.

822    17.    Wasserstein RL, Lazar NA. The ASA Statement on $p$ -Values: Context, Process, and

823          Purpose. Am Stat. 2016 Apr 2;70(2):129–33.

824    18.    Ioannidis JPA. Why most published research findings are false. PLoS Med. 2005 Aug

825          30;2(8):e124.

826    19.    Gundersen OE, Kjensmo S. State of the art: reproducibility in artificial intelligence.

827          AAAI. 2018 Apr 25;32(1).

828    20.    Verstynen T, Kording KP. Overfitting to 'predict' suicidal ideation. Nat Hum Behav.

829          2023 Apr 6;

830    21.    Riley P. Three pitfalls to avoid in machine learning. Nature. 2019 Aug;572(7767):27–9.

831    22.    Whalen S, Schreiber J, Noble WS, Pollard KS. Navigating the pitfalls of applying

832            machine learning in genomics. Nat Rev Genet. 2022 Mar;23(3):169–81.

833    23.    Vabalas A, Gowen E, Poliakoff E, Casson AJ. Machine learning algorithm validation

834            with a limited sample size. PLoS ONE. 2019 Nov 7;14(11):e0224365.

835    24.    Varma S, Simon R. Bias in error estimation when using cross-validation for model

836            selection. BMC Bioinformatics. 2006 Feb 23;7:91.

837    25.    Santos MS, Soares JP, Abreu PH, Araujo H, Santos J. Cross-Validation for

838            Imbalanced Datasets: Avoiding Overoptimistic and Overfitting Approaches [Research

839            Frontier]. IEEE Comput Intell Mag. 2018 Nov;13(4):59–76.

840    26.    Berisha V, Krantsevich C, Hahn PR, Hahn S, Dasarathy G, Turaga P, et al. Digital

841            medicine and the curse of dimensionality. npj Digital Med. 2021 Oct 28;4(1):153.

842    27.    Lones MA. How to avoid machine learning pitfalls: a guide for academic researchers.

843            arXiv. 2021;

844    28.    Kaufman S, Rosset S, Perlich C, Stitelman O. Leakage in data mining: Formulation,

845            detection, and avoidance. ACM Trans Knowl Discov Data. 2012 Dec 1;6(4):1–21.

846    29.    Rosenblatt M, Tejavibulya L, Jiang R, Noble S, Scheinost D. Data leakage inflates

847            prediction performance in connectome-based machine learning models. Nat Commun.

848            2024 Feb 28;15(1):1829.

849    30.    Bernett J, Blumenthal DB, Grimm DG, Haselbeck F, Joeres R, Kalinina OV, et al.

850            Guiding questions to avoid data leakage in biological machine learning applications.

851            Nat Methods. 2024 Aug 9;21(8):1444–53.

852    31.    Just MA, Pan L, Cherkassky VL, McMakin DL, Cha C, Nock MK, et al. Machine

853            learning of neural representations of suicide and emotion concepts identifies suicidal

854            youth. Nat Hum Behav. 2017 Oct 30;1:911–9.

855 32. Dukart J, Weis S, Genon S, Eickhoff SB. Towards increasing the clinical applicability
856     of machine learning biomarkers in psychiatry. Nat Hum Behav. 2021 Apr 5;5(4):431–2.

857 33. Kohavi R. A study of cross-validation and bootstrap for accuracy estimation and model
858     selection. International Joint Conference on Arti cial Intelligence. 1995;

859 34. Arlot S, Celisse A. A survey of cross-validation procedures for model selection. Stat
860     Surv. 2010;4(0):40–79.

861 35. Geisser S. The Predictive Sample Reuse Method with Applications. J Am Stat Assoc.
862     1975 Jun;70(350):320–8.

863 36. Bates S, Trevor H, Tibshirani R. Cross-Validation: What Does It Estimate and How
864     Well Does It Do It? Journal of the American Statistical Association.
865     2023;119(546):1434–45.

866 37. Krstajic D, Buturovic LJ, Leahy DE, Thomas S. Cross-validation pitfalls when selecting
867     and assessing regression and classification models. J Cheminform. 2014 Mar
868     29;6(1):10.

869 38. Hastie T, Tibshirani R, Friedman J. The Elements of Statistical Learning. 2nd ed. New
870     York, NY: Springer New York; 2009.

871 39. Yang L, Shami A. On hyperparameter optimization of machine learning algorithms:
872     theory and practice. Neurocomputing. 2020 Jul;

873 40. James G, Witten D, Hastie T, Tibshirani R. An Introduction to Statistical Learning. New
874     York, NY: Springer New York; 2013.

875 41. Bishop CM. Pattern recognition and machine learning. Springer New York; 2006.

876 42. Varoquaux G, Raamana PR, Engemann DA, Hoyos-Idrobo A, Schwartz Y, Thirion B.
877     Assessing and tuning brain decoders: Cross-validation, caveats, and guidelines.
878     Neuroimage. 2017 Jan 15;145(Pt B):166–79.

879 43. Martinez-Plumed F, Contreras-Ochando L, Ferri C, Hernandez-Orallo J, Kull M,

880    Lachiche N, et al. CRISP-DM Twenty Years Later: From Data Mining Processes to

881    Data Science Trajectories. IEEE Trans Knowl Data Eng. 2021 Aug 1;33(8):3048–61.

882    44.    Wirth R. CRISP-DM: Towards a standard process model for data mining. Proceedings

883    of the 4th international conference on the practical applications of knowledge

884    discovery and data mining; 2000.

885    45.    Chakraborty J, Majumder S, Menzies T. Bias in machine learning software: why?

886    how? what to do? Proceedings of the 29th ACM Joint Meeting on European Software

887    Engineering Conference and Symposium on the Foundations of Software Engineering.

888    New York, NY, USA: ACM; 2021. p. 429–40.

889    46.    Liang W, Tadesse GA, Ho D, Li F-F, Zaharia M, Zhang C, et al. Advances, challenges

890    and opportunities in creating data for trustworthy AI. Nat Mach Intell. 2022 Aug 17;

891    47.    Demšar J, Zupan B. Hands-on training about overfitting. PLoS Comput Biol. 2021 Mar

892    4;17(3):e1008671.

893    48.    Van Essen DC, Smith SM, Barch DM, Behrens TEJ, Yacoub E, Ugurbil K, et al. The

894    WU-Minn Human Connectome Project: an overview. Neuroimage. 2013 Oct 15;80:62–

895    79.

896    49.    Finn ES, Shen X, Scheinost D, Rosenberg MD, Huang J, Chun MM, et al. Functional

897    connectome fingerprinting: identifying individuals using patterns of brain connectivity.

898    Nat Neurosci. 2015 Nov;18(11):1664–71.

899    50.    He T, Kong R, Holmes AJ, Nguyen M, Sabuncu MR, Eickhoff SB, et al. Deep neural

900    networks and kernel regression achieve comparable accuracies for functional

901    connectivity prediction of behavior and demographics. Neuroimage. 2020 Feb

902    1;206:116276.

903    51.    Demeter DV, Engelhardt LE, Mallett R, Gordon EM, Nugiel T, Harden KP, et al.

904    Functional Connectivity Fingerprints at Rest Are Similar across Youths and Adults and

905    Vary with Genetic Similarity. iScience. 2020 Jan 24;23(1):100801.

906    52.    Schreiber J, Singh R, Bilmes J, Noble WS. A pitfall for machine learning methods
907           aiming to predict across cell types. Genome Biol. 2020 Nov 19;21(1):282.

908    53.    Yagis E, Atnafu SW, García Seco de Herrera A, Marzi C, Scheda R, Giannelli M, et al.
909           Effect of data leakage in brain MRI classification using 2D convolutional neural
910           networks. Sci Rep. 2021 Nov 19;11(1):22544.

911    54.    Cerqueira V, Torgo L, Mozetič I. Evaluating time series forecasting models: an
912           empirical study on performance estimation methods. Mach Learn. 2020
913           Nov;109(11):1997–2028.

914    55.    De Gooijer JG, Hyndman RJ. 25 years of time series forecasting. Int J Forecast. 2006
915           Jan;22(3):443–73.

916    56.    Samala RK, Chan H-P, Hadjiiski L, Helvie MA. Risks of feature leakage and sample
917           size dependencies in deep feature extraction for breast mass classification. Med Phys.
918           2021 Jun;48(6):2827–37.

919    57.    Emmanuel T, Maupong T, Mpoeleng D, Semong T, Mphago B, Tabona O. A survey on
920           missing data in machine learning. J Big Data. 2021 Oct 27;8(1):140.

921    58.    Poldrack RA, Huckins G, Varoquaux G. Establishment of best practices for evidence
922           for prediction: A review. JAMA Psychiatry. 2020 May 1;77(5):534–40.

923    59.    Vanwinckelen G, Blockeel H. Look before you leap: Some insights into learner
924           evaluation with cross-validation. Statistically Sound Data Mining. 2015 Nov 27;3–20.

925    60.    Badgeley MA, Zech JR, Oakden-Rayner L, Glicksberg BS, Liu M, Gale W, et al. Deep
926           learning predicts hip fracture using confounding patient and healthcare variables. npj
927           Digital Med. 2019 Apr 30;2:31.

928    61.    Dagaev N, Roads BD, Luo X, Barry DN, Patil KR, Love BC. A Too-Good-to-be-True
929           Prior to Reduce Shortcut Reliance. Pattern Recognit Lett. 2022 Dec;

930    62.    Hu F, Chen AA, Horng H, Bashyam V, Davatzikos C, Alexander-Bloch A, et al. Image

931       harmonization: A review of statistical and deep learning methods for removing batch

932       effects and evaluation metrics for effective harmonization. Neuroimage. 2023 Jul

933       1;274:120125.

934  63.  Bender R, Lange S. Adjusting for multiple testing--when and how? J Clin Epidemiol.

935       2001 Apr;54(4):343–9.

936  64.  García-Pérez MA. Use and misuse of corrections for multiple testing. Methods in

937       Psychology. 2023 Nov;8:100120.

938  65.  Thompson WH, Wright J, Bissett PG, Poldrack RA. Dataset decay and the problem of

939       sequential analyses on open datasets. eLife. 2020 May 19;9.

940  66.  Dwork C, Feldman V, Hardt M, Pitassi T, Reingold O, Roth A. The reusable holdout:

941       Preserving validity in adaptive data analysis. Science. 2015 Aug 7;349(6248):636–8.

942  67.  Hardt M, Ullman J. Preventing false discovery in interactive data analysis is hard. 2014

943       IEEE 55th Annual Symposium on Foundations of Computer Science. IEEE; 2014. p.

944       454–63.

945  68.  Chen C, Cao X, Tian L. Partial Least Squares Regression Performs Well in MRI-

946       Based Individualized Estimations. Front Neurosci. 2019 Nov 27;13:1282.

947  69.  Komeyer V, Eickhoff SB, Grefkes C, Patil KR, Raimondo F. A framework for

948       confounder considerations in AI-driven precision medicine. medRxiv. 2024 Feb 4;

949  70.  More S, Antonopoulos G, Hoffstaedter F, Caspers J, Eickhoff SB, Patil KR, et al.

950       Brain-age prediction: A systematic comparison of machine learning workflows.

951       Neuroimage. 2023 Apr 15;270:119947.

952  71.  Pourhoseingholi MA, Baghestani AR, Vahedi M. How to control confounding effects by

953       statistical analysis. Gastroenterol Hepatol Bed Bench. 2012;5(2):79–83.

954  72.  Snoek L, Miletić S, Scholte HS. How to control for confounds in decoding analyses of

955       neuroimaging data. Neuroimage. 2019 Jan 1;184:741–60.

956    73.    More S, Eickhoff SB, Caspers J, Patil KR. Confound removal and normalization in

957          practice: A neuroimaging based sex prediction case study. In: Dong Y, Ifrim G,

958          Mladenić D, Saunders C, Van Hoecke S, editors. Machine learning and knowledge

959          discovery in databases applied data science and demo track: european conference,

960          ECML PKDD 2020, ghent, belgium, september 14–18, 2020, proceedings, part V.

961          Cham: Springer International Publishing; 2021. p. 3–18.

962    74.    Hamdan S, Love BC, von Polier GG, Weis S, Schwender H, Eickhoff SB, et al.

963          Confound-leakage: Confound Removal in Machine Learning Leads to Leakage. arXiv.

964          2022;

965    75.    Mitchell M, Wu S, Zaldivar A, Barnes P, Vasserman L, Hutchinson B, et al. Model

966          cards for model reporting. Proceedings of the Conference on Fairness, Accountability,

967          and Transparency  - FAT* '19. New York, New York, USA: ACM Press; 2019. p. 220–

968          9.

969    76.    Walsh I, Fishman D, Garcia-Gasulla D, Titma T, Pollastri G, ELIXIR Machine Learning

970          Focus Group, et al. DOME: recommendations for supervised machine learning

971          validation in biology. Nat Methods. 2021 Oct;18(10):1122–7.

972    77.    Luo W, Phung D, Tran T, Gupta S, Rana S, Karmakar C, et al. Guidelines for

973          developing and reporting machine learning predictive models in biomedical research:

974          A multidisciplinary view. J Med Internet Res. 2016 Dec 16;18(12):e323.

975    78.    Norgeot B, Quer G, Beaulieu-Jones BK, Torkamani A, Dias R, Gianfrancesco M, et al.

976          Minimum information about clinical artificial intelligence modeling: the MI-CLAIM

977          checklist. Nat Med. 2020 Sep;26(9):1320–4.

978    79.    Gebru T, Morgenstern J, Vecchione B, Vaughan JW, Wallach H, Iii HD, et al.

979          Datasheets for datasets. Commun ACM. 2021 Dec;64(12):86–92.

980    80.    Holland S, Hosny A, Newman S, Joseph J, Chmielinski K. The Dataset Nutrition Label:

981          A Framework To Drive Higher Data Quality Standards. arXiv. 2018;

982    81.    Schwab S, Held L. Statistical programming: Small mistakes, big impacts. Significance.

983         2021 Jun;18(3):6–7.

984    82.    Barnes N. Publish your computer code: it is good enough. Nature. 2010 Oct

985         14;467(7317):753.

986    83.    Soares C. Is the UCI repository useful for data mining? Portuguese Conference on

987         Artificial Intelligence. 2003;209–23.

988    84.    van Giffen B, Herhausen D, Fahse T. Overcoming the pitfalls and perils of algorithms:

989         A classification of machine learning biases and mitigation methods. J Bus Res. 2022

990         May;144:93–106.

991    85.    Paleyes A, Urma R-G, Lawrence ND. Challenges in deploying machine learning: a

992         survey of case studies. ACM Comput Surv. 2022 Apr 30;