

Supplementary Materials for
**REFORMS: Consensus-based Recommendations for
Machine-learning-based Science**

Sayash Kapoor *et al.*

Corresponding author: Sayash Kapoor, sayashk@princeton.edu

Sci. Adv. **10**, eadk3452 (2024)
DOI: 10.1126/sciadv.adk3452

This PDF file includes:

Texts S1 and S2
References

Text S1: Guidelines for Filling Out the REFORMS Checklist

Visit `reforms.cs.princeton.edu` for the latest version.

These guidelines provide documentation for each item in the Consensus-based Reporting standards for ML-based science. We elaborate on why researchers should consider reporting the item, link to additional helpful resources to accomplish each item and add references to articles that describe the issues in depth.

We also provide a sample checklist based on Obermeyer et al. (102) (URL: <https://reforms.cs.princeton.edu/obermeyer-sample.pdf>).

As noted in our paper, some of the items in our checklist could be hard to report for specific studies. For instance, including a reproduction script to computationally reproduce all results (2e.) might not be possible for studies performed on academic computing clusters or those which use private data that cannot be released.

Instead of requiring strict adherence for each item, we suggest authors and referees decide which items are relevant for a study and how details can be reported better.

Module 1: Study design

The items in this section help communicate the purpose and goals of the study and how various decisions in the study design were arrived at. Details about the design of the study are important to clarify the applicability of the scientific claims of the study. They also help communicate the motivation behind researchers' various degrees of freedom, i.e., decisions researchers make throughout the research and analysis process that influence their findings.

1a. State the population or distribution about which the scientific claim is made.

Researchers make scientific claims about a given distribution or population that they are interested in studying. Note that this is the population of interest, and not the sample, which can be

specified later in (3b.)

To communicate the applicability of the claims, explicitly report the distribution or population about which you expect the scientific claims to hold. For example, “US children aged between 12 and 18” or “people engaging in online debates on climate change.”

1b. Describe the motivation for choosing this population or distribution (1a).

Justify why the researchers chose this population or distribution. For example: “We aimed to determine whether existing vaccines for COVID-19 are effective in children aged between 12 and 18. There are no prior studies on vaccine efficacy in this population.”

A valid motivation is having access to a dataset that inspired a research question, and thus the population or distribution of interest is limited by the dataset. For example, studying CDC data for all U.S. counties would limit the population of interest to US counties.

1c. Describe the motivation for the use of ML methods in the study.

Report the reasons for using ML methods and consider comparing it with alternative or traditional methods that could be used for similar aims.

For example, if the goal of the research is to make a prediction, i.e., if explanation is not a goal of the study, ML methods can help improve predictive accuracy.

See Hofman et al. (48) for an overview of the different types of modeling and their aims.

Module 2: Computational reproducibility

Computational reproducibility refers to the ability of a researcher to get the same figures and results that are reported in a paper or manuscript without making any changes to the code, data, or computing environment. This is important for ensuring the scientific validity of a study: errors can be uncovered quickly, independent researchers can verify the findings in a study, and researchers can easily build on a study’s results. Several journals currently require

computational reproducibility and have specific guidelines. If you're already using a discipline or journal-specific checklist, specify that here.

See Liu and Salganik (25) for a discussion on the importance and challenges of ensuring computational reproducibility.

Sandve et al. (106) discuss high-level imperatives and research practices that can enable computational reproducibility.

See the Social Science Data Editors' guidance on computational reproducibility (193).

Include as many of the items below as possible, in supplementary documents alongside a paper or pre-print that describes the study. Ideally, upload them to an established repository that provides a persistent identifier for the resources (such as Harvard Dataverse or Zenodo). Since code, data, and computational environments can have different versions over time, include the precise version that you use to generate the results reported in a study.

For some domains, sharing the code and dataset is not possible due to the presence of sensitive data. Specify below if such a restriction applies.

2a. Describe the dataset used for training and evaluating the model and provide a link or DOI to uniquely identify the dataset.

Report a permanent link or DOI to the specific version of the dataset used for training and evaluating the model. For a discussion of the importance of DOIs, see Peng et al. (97).

If an original dataset was used, also include the data dictionary for the dataset. A data dictionary describes metadata about the dataset, and familiarizes a reader to the properties and format of the data. The US Geological Survey has a detailed guide to data dictionaries, complete with examples and instructions (98).

If the dataset contains sensitive information and cannot be publicly released, consider releasing a synthetic dataset, or releasing the data per request or application. There are packages that support generation of a synthetic dataset such as synthpop for R (103).

2b. Provide details about the code used to train and evaluate the model and produce the results reported in the paper along with link or DOI to uniquely identify the version of the code used.

Provide a commit tag (for instance, on Github, GitLab, or BitBucket), a DOI, or equivalent documentation to precisely identify the version of the code used to train and evaluate the model and produce the exact results reported in the paper.

In the code, include comments with explanations of variables and operations to sufficiently mark different stages of the analysis for an unfamiliar reader. The documentation in (2d) can refer to these comments for greater clarity.

2c. Describe the computing infrastructure used.

To help readers understand the precise computing requirements for reproducing your study, whenever possible, report the following details on the infrastructure used to generate the results:

Hardware infrastructure: CPU, GPU, RAM, disk space. Operating system and its version.
Software environment: Programming language and version, documentation of all packages used along with versions and dependencies (e.g., through a requirements.txt file). An estimate of the time taken to generate the results.

Computing infrastructure is always changing, and thus could make it difficult or impossible to replicate a study with a slightly different environment. Having the exact details is crucial for replication.

See Requirements File Format (194) from Python's pip installer for an example of how to document package versions.

See Stodden and Miguez (109) for more detailed best practices to document computing infrastructure.

2d. Provide a README file which contains instructions for generating the results using the provided dataset and code.

Report the exact steps that should be taken by independent researchers to reproduce the results in your study, given access to the code, dataset, and computing environment specified in 2a-c.

A good README helps someone unfamiliar with the project by walking them through the steps of setting up and running the code provided, starting from environment requirements and installation, to examples of usage and expected results.

Consider using Nature’s README for software submission (195). See also the README template for social science replication packages (110).

The “Awesome README” repository compiles examples, templates, and best practices for writing README files (111).

2e. Provide a reproduction script to produce all results reported in the paper.

A script to produce all results reported in the paper using the code and dataset can substantially reduce the time it takes for an independent researcher to reproduce the results reported in a study.

The script should go through all steps involved in producing the results. For example, the script should download the packages, set the right dependencies, download and store the dataset in the correct location, set up the computational environment, pre-process the data, and run the code to produce exactly the same results as reported in the paper.

One option is a bash script which carries out each of the steps you list in (2d) (112). Another way is to use an online reproducibility platform such as CodeOcean, which allows researchers to share the required materials in 2a-c along with a reproduction script (113).

Note that this is a high bar for computational reproducibility, and in some cases, it might not be possible to provide such a script—for instance, if the analysis is run on an academic high-

performance computing cluster, or if the dataset does not allow for programmatic download. It could also be challenging to set up, and resources listed here might help. In case you are not able to share a reproduction script, specify why.

Comi (196) introduces CodeOcean for reproducible research, and shares how to create a CodeOcean capsule from Git.

Module 3: Data quality

This section is focused on reporting details about how the data used for developing and evaluating the model is collected. A good quality dataset is key to making valid scientific claims using ML models. The items in this section help readers understand and evaluate the quality of the data used in the modeling process.

3a. Describe source(s) of data, separately for the training and evaluation datasets (if applicable), along with the time when the dataset(s) are collected, the source and process of ground-truth annotations, and other data documentation.

Report details about the source of the dataset, separately for the training and validation data sets (if applicable). For instance, if re-using the dataset from a previous study, cite the study and explain what the source of the data collection was.

If collecting a new dataset, report the data collection process, who annotated the dataset, and how the annotations were carried out. Report the time-period and geographic locations of data collection.

You can also follow discipline-specific best-practices when releasing or using datasets. Examples include Datasheets for Datasets (99), Dataset Nutrition Labels (197), or the Brain Imaging Data Structure for Neuroimaging (198). If available, include such supplementary documents as supplementary materials along with the paper.

3b. State the distribution or set from which the dataset is sampled (i.e., the sampling frame).

The sampling frame is the source from which a sample is drawn (using a sampling method.)

The unit of the sampling frame is typically also the unit of the sample.

Report the sampling frame, which is the distribution or set from which the dataset is sampled. Include the sampling method (e.g., simple random, stratified, cluster sampling, etc.) Include any details about the distribution or population that pertains to the study (1a.).

Taherdoost (199) compiled a short guide to sampling in research.

3c. Justify why the dataset is useful for the modeling task at hand.

Report the rationale for why the dataset is useful for modeling and making the scientific claim reported in the study. Justifications could describe why the dataset is relevant to the modeling task, such as quantifying the population of interest well, or including novel insight that would be discovered through modeling.

3d. State the outcome variable of the model, along with descriptive statistics (split by class for a categorical outcome variable) and its definition.

The outcome or target variable of the ML model is the quantity that the model is used to predict, detect, classify, or estimate. In other words, it is the variable of interest in the modeling process.

Report the outcome variable of the ML model. Provide descriptive statistics (e.g., mean, median, and variance) for the outcome variable, if applicable. For tasks with a continuous outcome variable (i.e., regression tasks), consider providing a plot of the outcome's distribution, such as a histogram.

3e. State the sample size and outcome frequencies.

Report the total number of samples (for a tabular dataset, this is the total number of rows in the dataset) as well as the number of samples in each class for a classification task.

If there are individuals or entities with multiple observations, report both the number of distinct individuals, as well as overall rows or units of data. For example, if you have a dataset with 10,000 rows with data on 5,000 unique patients, report both of these numbers. See also (6b.)

3f. State the percentage of missing data, split by class for a categorical outcome variable.

Datasets often have missing samples. An estimate of missingness can give readers an idea of how important the methods for dealing with missing data are in a given study.

Report the number or percentage of missing samples for each feature, when possible. Alternatively, provide summary statistics for the proportion of missing data.

See also (4c.) for methods for handling missing data.

3g. Justify why the distribution or set from which the dataset is drawn (3b.) is representative of the one about which the scientific claim is being made (1a.).

Justify why the distribution or set from which the dataset is drawn (3b.) is representative of the population about which the scientific claim is being made (1a.).

There are many reasons the sampling frame could be unrepresentative: for example, if it is a convenience sample, if it under-represents minorities, or constitutes a too small sample size (13). If the sample is unrepresentative of the target population, note this as a concern in the section on external validity (8a. Evidence of external validity).

Module 4: Data preprocessing

Pre-processing is the series of steps taken to convert the dataset used from its raw form into the final form used in the modeling process. This includes data selection (i.e., selecting a set of samples from the dataset to be included in the modeling process) as well as other transformations of the data, such as imputing missing data and normalizing feature values.

Since pre-processing steps can influence the scientific claims made based on ML models (28), it is important to specify the exact steps used in a study.

4a. Describe whether any samples are excluded with a rationale for why they are excluded.

Researchers might exclude some samples from the dataset—for instance, to remove outliers or to only focus on certain subsets. Report the criteria for selecting a subset of rows from the initial dataset (if any).

4b. Describe how impossible or corrupt samples are dealt with.

Some datasets might have feature values that are impossible (for instance, if the height of a human is recorded as greater than 10 feet). Some samples might have corrupt data.

Report the checks made for impossible or corrupt data. In case you find impossible or corrupt data, report mitigation strategies, such as methods used for detecting or removing outliers.

4c. Describe all transformations of the dataset from its raw form (3a.) to the form used in the model, for instance, treatment of missing data and normalization—preferably through a flow chart.

Researchers often perform several transformations on a dataset before using it in an ML model. For example, they might impute missing data in a dataset using mean imputation or over-sample data from the minority class.

Report the precise sequence of all transformations of data from its raw form to the final

form used in the model (e.g., missing data imputation, feature or outcome normalization, data augmentation using oversampling), preferably through a flow-chart, like a STROBE flow diagram (200).

Specify if each transformation is data-dependent (e.g., mean imputation) or data-independent (e.g., log transformation). Note that data-dependent transformations must be done within splits. For example, when using 5-fold cross-validation, perform mean imputation within each of the folds instead of performing it on the entire data together to avoid leaking information between the training and test data. See also 6a.

(136) discuss how poorly imputed data can lead to poor interpretability of the final model.

Module 5: Modeling

There are many steps involved in creating an ML model. This makes it hard to report the exact details of how an ML model is created, and can hinder replication by independent researchers. Specify the main steps in the modeling process, including feature selection, the types of models considered, and evaluation.

5a. Describe, in detail, all models trained.

To help readers determine how the models were trained, provide a detailed description of all models trained over the course of the study. For each model, include: Inputs (including any feature selection steps and a description of the set of features used) and outputs Types of models implemented (e.g., Random Forests, Neural Networks) Loss function used

5b. Justify the choice of model types implemented.

Describe why the types of models used are relevant for the study. Examples are “using a standard method for this field such as regularized regressions”, or “using decision trees for high explainability.”

(8) describe various ML models that are suitable for different modeling tasks.

5c. Describe the method for evaluating the model(s) reported in the paper, including details of train-test splits or cross-validation folds.

Evaluating ML models requires testing them on data that they were not trained on, for instance by using a held-out test set or cross-validation (CV).

Report how the dataset is split for evaluating the ML model(s), for instance: Cross-validation or nested CV Held-out test set (internal validation set) True out-of-sample set (external validation set; where the data comes from a different set compared to training data)

For the model evaluation method used, report details such as the number of samples in each train-test split or CV fold, as well as the number of samples of each class in each split (for a classification task).

Documentation from the Python package scikit learn elaborates why and how to do a train-validation-test split (201).

Vehtari (202) describes various scenarios where using CV is appropriate.

Neunhoeffler and Sternberg (79) highlight a common failure mode: using CV for both model selection and evaluation. Using nested CV helps address this issue.

Cawley and Talbot (144) explore this issue in more detail and describe procedures for nested CV.

5d. Describe the method for selecting the model(s) reported in the paper.

Several ML models might be fit using the training set.

Report the criteria for choosing the final model(s) reported in the study. For instance, report if model performance on the training set, internal cross-validation fold (for nested cross-validation) or a separate validation set was used to select the final model(s) reported in the paper.

Raschka (*143*) gives an overview of model selection techniques.

5e. For the model(s) reported in the paper, specify details about the hyperparameter tuning.

ML models often have hyperparameters. For example, Lasso regression has an additional penalty term (lambda or λ) that can be tuned. Tuning hyperparameters—trying different values and picking the one that works best—can help find the optimal performance for a given model and dataset.

Report the method used to compare the performance of different hyperparameter values. This should include details of what values for each parameter are considered, why these values are reasonable, how various hyperparameters are selected (for example, grid search or random search (*201*)), and which hyperparameters are used in the final model(s) reported in the paper.

5f. Justify that model comparisons are against appropriate baselines.

If comparing model performance against baselines, justify how the baselines are tuned appropriately and the model comparison is fair if applicable. (Note that this does not apply to comparisons against non-model based performance, such as comparing ML methods with human performance.)

Sculley et al. (*24*) highlight several results in ML research that compare against weak baselines.

Lin (*155*) highlights that comparisons against weak baselines can result in overoptimism.

Module 6: Data leakage

Data leakage is a spurious relationship between the independent variables and the target variable that arises as an artifact of the data collection, sampling, pre-processing or modeling steps. Since the spurious relationship won't be present in the distribution about which scientific claims

are made, leakage usually leads to inflated estimates of model performance. Items in this section help detect and prevent leakage in the models developed and evaluated in a study.

Kapoor and Narayanan (39) discuss the prevalence of leakage and provide “Model Info Sheets” to detect and prevent leakage in ML-based science.

6a. Justify that pre-processing (Module 4) and modeling (Module 5) steps only use information from the training dataset (and not the test dataset).

When information from the test set is used during the training process, it leads to overly optimistic performance and results in data leakage.

Justify how all pre-processing (Module 4) and modeling (Module 5) steps only use information from the training data and not the entire dataset (e.g., they were performed after the data splits or cross-validation splits).

Vandewiele et al. (81) show how oversampling before partitioning the training data and test data can cause errors in models, with several studies incorrectly reporting near-perfect accuracy.

6b. Describe methods used to address dependencies or duplicates between the training and test datasets (e.g. different samples from the same patients are kept in the same dataset partition).

In some cases, samples in the dataset might have dependencies. For example, a clinical dataset might have many samples from the same patient. In such cases, the train-test split or cross-validation (CV) split should take these dependencies into account—for instance, by including all samples from each patient in the same CV fold or train-test split.

Similarly, duplicates in the datasets can also spread across training and test sets if the dataset is split randomly. This should be avoided, as it leaks information across the train-test split.

Report if the dataset used has dependencies or duplicates. If it does, detail how these are addressed (for example, by using block CV or removing duplicate rows of data).

Malik (162) outlines alternatives for CV that helps reduce dependencies.

Bergmeir and Benítez (*164*) find that blocked CV for time series evaluation deals with temporal autocorrelation.

Hammerla and Plotz (*165*) demonstrate how neighborhood bias can affect data recordings close in time and introduce “meta-segmented CV” to deal with such dependencies.

Roberts et al. (*166*) describe block CV strategies for a number of structures with dependencies, including temporal, spatial, and hierarchical dependencies.

6c. Justify that each feature or input used in the model is legitimate for the task at hand and does not lead to leakage.

Leakage can result from any of the features used in a model being a proxy for the outcome. For example, Filho et al. (*167*) found that a prominent paper on hypertension prediction (*168*) suffered from data leakage due to illegitimate features. The model included the use of anti-hypertensive drugs as a feature in a clinical model used to predict hypertension.

Justify why each of the features used in the model is legitimate for the task at hand. Note that you do not necessarily need to list each feature individually; instead, you can provide arguments for a set of features together in case the same argument applies to all of them.

Module 7: Metrics and uncertainty

The performance of ML models is key to the scientific claims of interest. Since there are many possible choices that authors can make when choosing performance metrics, it is important to reason about why the metrics used are appropriate for the task at hand. Additionally, communicating and reasoning about uncertainty is important to discourage readers from ignoring the uncertainty in the final results.

7a. State all metrics used to assess and compare model performance (e.g., accuracy, AU-ROC etc.). Justify that the metric used to select the final model is suitable for the task.

Several metrics are often used to assess how well an ML model performs and to compare the performance of different ML models. In some cases, these metrics are reported as part of a paper's final results, while in others, they are used to make intermediate decisions such as identifying which models to include in the study or to decide which hyperparameters should be used.

Report all metrics used to assess and compare model performance (e.g., Accuracy, AUC-ROC etc.). Include metrics that are used to make decisions about which model(s) are reported as well as the metrics used to evaluate the reported model(s).

Some metrics are unsuitable for certain problems. For example, accuracy might not be suitable to measure the performance of an ML model in the presence of heavy class imbalance (see (8), Table 4). Justify the choice of metric(s) used for the scientific claim being made based on the ML model.

7b. State uncertainty estimates (e.g., confidence intervals, standard deviations), and give details of how these are calculated.

For each performance metric reported in a paper, report an estimate of uncertainty such as standard deviations or confidence intervals. This could be part of graphs or tables in the paper.

Note that applying a bootstrap on the validation set is one way to get uncertainty estimates for a population mean based on a sample from that population.

Report the uncertainty estimate. Also report how the uncertainty estimate is calculated and justify why the method used for uncertainty estimation is valid.

Simmonds et al. (38) outline the different sources of uncertainty that should be quantified in a study.

Raschka (143) walks through bootstrapping to obtain an uncertainty estimate.

7c. Justify the choice of statistical tests (if used) and a check for the assumptions of the statistical test.

Statistical tests used for comparing model performance come with several assumptions.

Report the type of statistical test used in the paper (if any) for comparing model performance. Report the assumptions of the statistical test and justify why these assumptions are satisfied.

If using bootstrapped confidence intervals for performance metrics, one statistical test is to see if the interval contains a baseline value. Raschka (*143*) outlines various statistical tests for comparing supervised learning algorithms. Note that reliance on statistical significance testing has led to misinterpretations and false conclusions (*179*).

Module 8: Generalizability and limitations

8a. Describe evidence of external validity.

External validity (or “generalizability”) refers to the applicability of a scientific claim beyond the specific dataset based on which it is made. This includes the extent to which the findings from a study’s sample apply to the target population, as well as the extent to which the findings apply to other populations, outcomes, and contexts (*181*). For example, evaluating an ML model on a different dataset or a new clinical setting that it was not trained on is a test of its external validity.

Researchers can use a mix of quantitative and theoretical approaches to make arguments regarding their findings’ ability to generalize to other populations, outcomes, and contexts. They can report quantitative evidence by testing their claims in out-of-distribution data. They can make theoretical arguments about their expectations of external validity by referring to prior literature and reasoning about the level of similarity between contexts (*63*).

Report evidence regarding the external validity of the study’s findings.

8b. Describe contexts in which the authors do not expect the study’s findings to hold.

Explicit boundaries around the applicability of a scientific claim can help clarify which settings we should expect the scientific claims to hold in. Authors are in the best position to understand limits to the applicability of their claims.

Report examples of settings or domains where the scientific claims made in the study do not hold.

Raji et al. (182) discuss issues with ML models used in real-world settings. These issues stem in part from a lack of focus on identifying when models are not expected to work.

Text S2: Table of References on Reporting Quality & Problems in Scientific Literature

This appendix provides additional details on some of the citations from the main text. We include references from the main text that address: (1) the quality of reporting in past scientific literature, or (2) examples of problems that occurred in past scientific literature. This appendix does not constitute a comprehensive list of all published references on these topics. The table has 44 entries with details about their relevance to our review.

The citations are listed in order of appearance in the main text, with section headings corresponding to the headings from the text. Some sections from the main text are omitted because they do not contain references that match our criteria for inclusion in the table. Some citations are included in the table more than once because they appear in multiple sections. Many of the references focus specifically on machine learning (ML)-based science, but we also include references about science with traditional statistical methods because some of the best practices and shortcomings are shared in ML-based science and other quantitative sciences.

Reference	Findings about reporting quality in past literature or problems in past literature	Discipline	Literature examined	ML-Focused?
MODULE 1: STUDY GOALS				
Introduction				
Hofman et al., 2017, "Prediction and explanation in social systems" (28)	The authors re-evaluate data from a prior paper to demonstrate how different (but equally reasonable) choices in research design can lead to different results from the same data. This includes an example of how slight differences in the definition of a research question can lead to substantially different results.	Computational social science	Re-evaluation of data from 1 prior paper on prediction of information cascade size on Twitter	Yes
1a) Population or distribution about which the scientific claim is made				
Lundberg et al., 2021, "What Is Your Estimand? Defining the Target Quantity Connects Statistical Evidence to Theory" (58)	Only 9 out of 32 papers (28%) provided sufficient information for a reader to "confidently" identify the target population about which the scientific claim is made (p. 553).	Sociology	32 quantitative papers in 2018 volume of a top sociology journal	No
Tooth et al., 2005, "Quality of Reporting of Observational Longitudinal Research" (62)	33 out of 49 papers (67%) define a target population.	Epidemiology & medicine	49 longitudinal studies on strokes in six journals, 1999-2003	No
MODULE 2: COMPUTATIONAL REPRODUCIBILITY				
Introduction				
Verstynen and Kording, 2023, "Overfitting to 'predict' suicidal ideation" (83)	The code for the feature selection step in a flawed prior paper was not released, so Verstynen and Kording could not pinpoint the exact source of errors.	Psychology, neuroscience, and biomedical engineering	1 paper on prediction of suicidal ideation	Yes
Current computational reproducibility standards fall short				
Stodden et al., 2018, "An empirical analysis of journal policy effectiveness for computational reproducibility" (85)	Stodden et al. attempted to contact the authors of 204 papers published in the journal Science to obtain reproducibility materials. Only 44% of authors responded.	Multi-disciplinary	204 quantitative papers in Science	No
Gabelica et al., 2022, "Many researchers were not compliant with their published data sharing statement: A mixed-methods study" (86)	Gabelica et al. examined 333 open-access journals indexed on BioMed Central in January 2019 and found that out of the 1,792 papers that pledged to share data upon request, 1,669 did not do so, resulting in a 93% data unavailability rate.	Biology, health sciences and medicine	1,792 papers published in 333 BioMed Central open-access journals in January 2019	No
Vasilevsky et al., 2017, "Reproducible and reusable research: Are journal data sharing policies meeting the mark?" (87)	Vasilevsky et al. examined the data-sharing policies of 318 biomedical journals and discovered that almost one-third lacked any such policies, and those that did often lacked clear guidelines for author compliance.	Biology, health sciences and medicine	318 biomedical journals (Biochemistry and Molecular Biology, Cell Biology, Crystallography, Developmental Biology, Biomedical Engineering, Immunology, Medical Informatics, Microbiology, Microscopy, Multidisciplinary Sciences, and Neurosciences)	No
Computational reproducibility allows independent researchers to find errors in original papers				

Hofman et al., 2021, "Expanding the scope of reproducibility research through data analysis replications" (80)	Hofman et al. analyze 11 papers and find various shortcomings in this body of literature.	Multi-disciplinary	11 computational social science papers	No
Vandewiele et al., 2021, "Overly optimistic prediction results on imbalanced data: A case study of flaws and benefits when applying over-sampling" (81)	Vandewiele et al. analyze 24 papers on pre-term birth prediction and find 21 of these papers suffer from leakage.	Medicine	24 papers on pre-term risk prediction	Yes
MODULE 3: DATA QUALITY				
3a) Data source(s)				
Navarro et al., 2022, "Completeness of reporting of clinical prediction models developed using supervised machine learning: a systematic review" (114)	98% of articles adhered to the guidelines for reporting data source from the TRIPOD statement.	Epidemiology & medicine	152 articles on diagnostic or prognostic prediction models across medical fields, published 2018-2019	Yes
Yusuf et al., 2020, "Reporting quality of studies using machine learning models for medical diagnosis: a systematic review" (115)	24 out of 28 papers (86%) reported information about their data source, defined as "Where and when potentially eligible participants were identified (setting, location and dates)" (p. 3).	Medicine	28 "medical research studies that used ML methods to aid clinical diagnosis," published July 2015-July 2018	Yes
Kim et al., 2016, "Garbage in, Garbage Out: Data Collection, Quality Assessment and Reporting Standards for Social Media Data Use in Health Research, Infodemiology and Digital Disease Detection" (116)	Studies that utilize social media data frequently omit important information about their data collection process, such as details about the development and assessment of search filters. This paper provides a framework for reporting this information.	Health media	Studies that use social media data (this is not a formal review paper, but it provides several examples)	No
Geiger et al., 2020, "Garbage In, Garbage Out? Do Machine Learning Application Papers in Social Computing Report Where Human-Labeled Training Data Comes From?" (117)	There was "wide divergence" in whether papers followed best practices for reporting the data annotation process, such as reporting: "who the labelers were, what their qualifications were, whether they independently labeled the same items, whether inter-rater reliability metrics were disclosed, what level of training and/or instructions were given to labelers, whether compensation for crowdworkers is disclosed, and if the training data is publicly available" (p. 325).	Multi-disciplinary: "the papers represented political science, public health, NLP, sentiment analysis, cybersecurity, content moderation, hate speech, information quality, demographic profiling, and more" (p. 328)	164 "machine learning application papers... that classified tweets from Twitter" (p. 326)	Yes
3b) Sampling frame				

Navarro et al., 2022, "Completeness of reporting of clinical prediction models developed using supervised machine learning: a systematic review" (114)	105 out of 152 studies (69%) reported their eligibility criteria.	Epidemiology & medicine	152 articles on diagnostic or prognostic prediction models across medical fields, published 2018-2019	Yes
Tooth et al., 2005, "Quality of Reporting of Observational Longitudinal Research" (62)	41 out of 49 papers (84%) reported their sampling frame, and 32 out of 49 papers (65%) reported their eligibility criteria.	Epidemiology & medicine	49 longitudinal studies on strokes in six journals, 1999-2003	No
Porzolt et al., 2019, "Inclusion and exclusion criteria and the problem of describing homogeneity of study populations in clinical trials" (118)	75 out of 100 studies (75%) reported inclusion criteria. 6 of those 75 studies (8%) also reported exclusion criteria.	Medicine	100 publications on "quality of life" assessments	No
3d) Outcome variable				
Credé and Harms, 2021, "Three cheers for descriptive statistics—and five more reasons why they matter" (122)	In a review of literature that was still a work-in-progress at the time Credé and Harms published this commentary, "Among the articles coded to date, less than half report the ethnicity of the participants or the types of jobs held by the participants and only 56% report data on the industry in which the data were collected. Other interesting—and to meta-analysts potentially important—information is also remarkably often unreported" (p. 486). (Note: This commentary discusses descriptive statistics broadly, not just descriptive statistics for outcome variables.)	Industrial and organizational psychology	Articles from four top journals in industrial and organizational psychology (number of articles is not reported)	No
Larson-Hall and Plonsky, 2015, "Reporting and interpreting quantitative research findings: What gets reported and recommendations for the field" (123)	Meta-analyses frequently had to omit large numbers of primary articles from their analyses due to insufficient descriptive statistics in the primary articles. (Note: This article discusses descriptive statistics broadly, not just descriptive statistics for outcome variables.)	Second language acquisition	Approximately 90 meta-analyses in second language acquisition	No
3e) Sample size				
Plonsky, 2013, "Study Quality in SLA: An Assessment of Designs, Analyses, and Reporting Practices in Quantitative L2 Research" (125)	99% of studies reported sample size.	Second language acquisition	606 studies in second language acquisition journals, published 1990-2010	No
Tooth et al., 2005, "Quality of Reporting of Observational Longitudinal Research" (62)	100% of 49 longitudinal studies reported the total number of participants from the first wave of their study. However, only 25 out of 49 (51%) reported the number of participants after attrition at each subsequent wave.	Epidemiology & medicine	49 longitudinal studies on strokes in six journals, 1999-2003	No
3f) Missingness				
McKnight et al., 2007, "Missing Data: A Gentle Introduction" (126)	Around 90% of articles had missing data, and the average amount of missing data per study was over 30%. Furthermore, "few of the articles included explicit mention of missing data, and even fewer indicated that the authors attended to missing data, either by performing statistical procedures or by making disclaimers regarding the studies in the results and conclusions" (p. 3).	Psychology	Over 300 publications from a prominent psychology journal	No

Peugh and Enders, 2004, "Missing Data in Educational Research: A Review of Reporting Practices and Suggestions for Improvement" (127)	Among the articles Peugh and Enders reviewed, "[d]etails concerning missing data were seldom reported" and "[t]he methods used to handle missing data were, in many cases, difficult to ascertain because explicit descriptions of missing-data procedures were rare" (p. 537). However, Peugh and Enders were able to infer the amount of missingness in some studies by examining the "discrepancy between the reported degrees of freedom for a given analysis and the degrees of freedom that one would expect on the basis of the stated sample size and design characteristics" (p. 537). In articles published in 1999, they detected missing data in 16% of studies, but they write that this is likely a "gross underestimate" of the actual prevalence of missing data. Among articles published in 2003, they were able to detect missing data in 42% of articles, which is higher than in 1999 due to changes in reporting practices following a recommendation by an American Psychological Association task force.	Educational research	989 studies published in 1999 and 545 studies published in 2003 in 23 applied educational research journals	No
Salganik et al., 2020, Supplementary information for "Measuring the predictability of life outcomes using a scientific mass collaboration" (34)	There are many reasons for missing data in survey data, including a respondent not participating in a given wave of a longitudinal survey, respondents refusing to answer some questions, skip patterns in the survey design, and redaction for privacy. In a modified version of a well-known, high-quality social survey dataset, 73% of possible data entries were missing, and the largest source of missingness was survey skip patterns. This high level of missingness emphasizes the importance of careful attention to handling missing data.	Sociology	1 study with a well-known social survey data set	Yes
Nijman et al., 2022, "Missing data is poorly handled and reported in prediction model studies using machine learning: a literature review" (129)	"A total of 56 (37%) prediction model studies did not report on missing data and could not be analyzed further. We included 96 (63%) studies which reported on the handling of missing data. Across the 96 studies, 46 (48%) did not include information on the amount or nature of the missing data" (p. 220).	Medicine	152 ML-based clinical prediction model studies, published 2018-2019	Yes
Navarro et al., 2022, "Completeness of reporting of clinical prediction models developed using supervised machine learning: a systematic review" (114)	"Forty-four studies reported how missing data were handled (28.9%, 95% CI 22.3 to 36.6). The missing data item consists of four sub-items of which three were rarely addressed in included studies. Within 28 studies that reported handling of missing data: three studies reported the software used (10.7%, CI 3.7 to 27.2), four studies reported the variables included in the procedure (14.3%, CI 5.7 to 31.5) and no study reported the number of imputations (0%, CI 0.0 to 39.0)" (pp. 6-7).	Epidemiology & medicine	152 articles on diagnostic or prognostic prediction models across medical fields, published 2018-2019	Yes
Little et al., 2013, "On the Joys of Missing Data" (130)	"Among the 80 reviewed studies, only 45 (56.25%) mentioned missing data explicitly in the text or a table of descriptive statistics. Of those 45, only three mentioned testing whether the missingness was related to other variables, justifying their [missingness at random] assumption" (p. 156).	Pediatric psychology	80 empirical studies in the 2012 issues of a pediatric psychology journal	No
Nicholson et al., 2016, "Attrition in developmental psychology" (131)	Among 541 longitudinal studies, only 253 (47%) discussed missingness due to attrition, and only 99 (18%) explicitly discussed whether missingness due to attrition was "missing at random," "missing completely at random," or "missing not at random."	Developmental psychology	541 longitudinal studies in major developmental journals, published 2009 and 2012	No

Sterner, 2011, "What Is Missing in Counseling Research? Reporting Missing Data" (132)	In the first journal, "14 of 66 (21%) articles referenced missing data on some level. Of these 14 articles, 11 mentioned missing data specifically... In the remaining 52 JCD articles, no information was provided on whether missing data existed." In the second journal, "one of 28 (4%) empirically based research articles made reference to screening for missing data; however, no mention was made of missing data in the remaining articles" (p. 56).	Counseling	94 empirical research articles in two top counseling journals, published 2004 to 2008	No
Tooth et al., 2005, "Quality of Reporting of Observational Longitudinal Research" (62)	Only 19 out of 49 articles (39%) reported on missing data items at each longitudinal wave, and only 2 out of 42 articles (5%) that had missing data in their analyses described imputation, weighting, or sensitivity analyses for handling missing data.	Epidemiology & medicine	49 longitudinal studies on strokes in six journals, 1999-2003	No
Hussain et al., 2017, "Quality of missing data reporting and handling in palliative care trials demonstrates that further development of the CONSORT statement is required: a systematic review" (133)	101 out of 108 studies (94%) reported the number of participants who were missing in the primary outcome analysis; however, reporting rates were lower for other details about missing data and for methods of handling missing data.	Epidemiology	108 articles on palliative care randomized controlled trials, published 2009-2014	No
3g) Dataset for evaluation is representative				
Tooth et al., 2005, "Quality of Reporting of Observational Longitudinal Research" (62)	Among several reporting criteria this review examined, "the criteria in the checklist representing selection bias were the least frequently reported overall" (p. 285). Specifically, selection-in biases were discussed in 14 out of 49 articles (28%), comparison of consenters with non-consenters was discussed in 1 out of 47 applicable articles (2%), and loss to follow-up was accounted for in the analyses of 1/41 applicable articles (5%). Additionally, 37 out of 49 articles (75%) discuss how their results relate to the target population.	Epidemiology & medicine	49 longitudinal studies on strokes in six journals, 1999-2003	No
MODULE 4: DATA PREPROCESSING				
4c) Data transformations				
Vandewiele et al., 2021, "Overly optimistic prediction results on imbalanced data: a case study of flaws and benefits when applying over-sampling" (81)	Vandewiele et al. analyze 24 papers on pre-term birth prediction and find 11 of these papers improperly transform data (by oversampling before splitting into train and test sets).	Medicine	24 papers on pre-term risk prediction	Yes
MODULE 5: MODELING				
5d) Model selection method				
Neunhoeffer and Sternberg, 2019, "How Cross-Validation Can Go Wrong and What to Do About It." (79)	Neunhoeffer and Sternberg demonstrate that the main findings of a prominent political science paper fail to reproduce due to improper model selection. In particular, model selection was done on the same data that was used for evaluation.	Political Science	1 prominent political science paper	Yes
5e) Hyper-parameter selection				
Dodge et al., 2019, "Show Your Work: Improved Reporting of Experimental Results" (149)	Dodge et al. find that among 50 random papers from a prominent natural language processing conference, while 74% of papers reported at least some information about the best performing hyperparameters, 10% of fewer reported more specific details about hyperparameter search or the effect of hyperparameters on performance.	Natural language processing	50 random papers from a prominent natural language processing conference in 2018	Yes
5f) Appropriate baselines				

Sculley et al., 2018, “Winner’s curse? On pace, progress, and empirical rigor” (24)	Sculley et al. discuss five papers that provide evidence of improper comparison with baselines in different areas of ML, suggesting that empirical progress in the field can be misleading.	ML	5 papers identifying poor performance compared to baselines in different areas of ML	Yes
MODULE 6: DATA LEAKAGE				
Introduction				
Kapoor and Narayanan, 2022, “Leakage and the reproducibility crisis in ML-based science” (39)	Kapoor and Narayanan found that leakage affects hundreds of papers across 17 fields.	Multi-disciplinary	A survey of leakage issues across 17 fields	Yes
Train-test separation is maintained				
Poldrack et al., 2020, “Establishment of best practices for evidence for prediction: A review” (160)	Poldrack et al. find that of the 100 neuropsychiatry studies that claimed to predict patient outcomes, 45 only reported in-sample statistical fit as evidence for predictive accuracy.	Neuropsychiatry	100 published studies between December 24, 2017 and October 30, 2018 in PubMed using search terms “fMRI prediction” and “fMRI predict”	Yes
Dependencies or duplicates between datasets				
Roberts et al., 2021, “Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans” (15)	Roberts et al. discuss the issue of “Frankenstein” datasets: datasets that combine multiple other sources of data and can end up using the same data twice—for instance, if two datasets rely on the same underlying data source are combined into a larger dataset.	Medicine	62 studies that claimed to diagnose or prognose Covid-19 using chest x-rays	Yes
MODULE 7: METRICS AND UNCERTAINTY				
7b) Uncertainty estimates				
Simmonds et al., 2022, “How is model-related uncertainty quantified and reported in different disciplines?” (38)	Simmonds et al. show that across seven fields, no fields consistently reported complete model uncertainties, and that the type of uncertainties reported varied by field.	Multi-disciplinary	496 studies across 7 fields that included statistical models	No
MODULE 8: GENERALIZABILITY AND LIMITATIONS				
Introduction				
Raji et al., 2022, “The Fallacy of AI Functionality” (182)	Raji et al. review real-world applications of technologies that claim to use ML and categorize several ways in which such technology frequently failed, including “lack of robustness to changing external conditions” (p. 9).	Computer science and law (real-world ML applications)	283 cases of failures of technology that claimed to be AI, ML or data-driven between 2012 to 2021	Yes
Liao et al., 2021, “Are We Learning Yet? A Meta-Review of Evaluation Failures Across Machine Learning” (183)	Liao et al. find that the same types of evaluation failures occur across a wide range of ML tasks and algorithms. They provide a taxonomy of common internal and external validity failures.	Computer science	107 “survey papers from computer vision, natural language processing, recommender systems, reinforcement learning, graph processing, metric learning, and more”	Yes

Reporting on external validity falls short in past literature				
Tooth et al., 2005, "Quality of Reporting of Observational Longitudinal Research" (62)	37 out of 49 papers (75%) discuss how the findings from their sample generalize to their target population, and 26 out of 49 papers (53%) discuss generalizability beyond the target population.	Epidemiology & medicine	49 longitudinal studies on strokes in six journals, 1999-2003	No
Bozkurt et al., 2020, "Reporting of demographic data and representativeness in machine learning models using electronic health records" (190)	The authors argue that descriptive statistics about the study sample should be provided in order to be transparent about representativeness of the target population. They find that of 164 studies that trained ML models with electronic health records data, "Race/ethnicity was not reported in 64%; gender and age were not reported in 24% and 21% of studies, respectively. Socioeconomic status of the population was not reported in 92% of studies." They also find, "Few models (12%) were validated using external populations" (p. 1878).	Medicine	164 studies that trained ML models with electronic health records data	Yes
Navarro et al., 2023, "Systematic review finds 'spin' practices and poor reporting standards in studies on machine learning-based prediction models" (191)	"In the main text, 86/152 (56.6% [95% CI 48.6 - 64.2]) studies made recommendations to use the model in clinical practice, however, 74/86 (86% [95% CI 77.2 - 91.8]) lacked external validation in the same article. Out of the 13/152 (8.6% [95% CI 5.1 - 14.1]) studies that recommended the use of the model in a different setting or population, 11/ 13 (84.6% [95% CI 57.8 - 95.7]) studies lacked external validation" (p. 104).	Epidemiology & medicine	152 articles on diagnostic or prognostic prediction models across medical fields, published 2018-2019	Yes

REFERENCES AND NOTES

1. S. Athey, G. W. Imbens, Machine learning methods that economists should know about. *Annu. Rev. Econ.* **11**, 685–725 (2019).
2. D. R. Schrider, A. D. Kern, Supervised machine learning for population genetics: A new paradigm. *Trends Genet.* **34**, 301–312 (2018).
3. J. J. Valletta, C. Torney, M. Kings, A. Thornton, J. Madden, Applications of machine learning in animal behaviour studies. *Anim. Behav.*, **124**, 203–220 (2017).
4. R. Iniesta, D. Stahl, P. McGuffin, Machine learning, statistical learning and the future of biological research in psychiatry. *Psychol. Med.* **46**, 2455–2465 (2016).
5. S. Tonidandel, E. B. King, J. M. Cortina, Big data methods: Leveraging modern data analytic techniques to build organizational science. *Organ. Res. Methods* **21**, 525–547 (2018).
6. T. Yarkoni, J. Westfall, Choosing prediction over explanation in psychology: Lessons from machine learning. *Perspect. Psychol. Sci.* **12**, 1100–1122 (2017).
7. J. Grimmer, M. E. Roberts, B. M. Stewart, Machine learning for social science: An agnostic approach. *Annu. Rev. Polit. Sci.* **24**, 395–419 (2021).
8. A. K. Leist, M. Klee, J. H. Kim, D. H. Rehkopf, S. P. A. Bordas, G. Muniz-Terrera, S. Wade, Mapping of machine learning approaches for description, prediction, and causal inference in the social and health sciences. *Sci. Adv.* **8**, eabk1942 (2022).
9. T. L. Wiemken, R. R. Kelley, Machine learning in epidemiology and health outcomes research, *Annu. Rev. Public Health* **41**, 21–36 (2020).
10. H. R. Varian, Big data: New tricks for econometrics. *J. Econ. Perspect.* **28**, 3–28 (2014).
11. S. Mullainathan, J. Spiess, Machine learning: An applied econometric approach. *J. Econ. Perspect.* **31**, 87–106 (2017).
12. L. Breiman, Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Stat. Sci.* **16**, 199–231 (2001).
13. J. Hullman, S. Kapoor, P. Nanayakkara, A. Gelman, A. Narayanan, “The worst of both worlds: A comparative analysis of errors in learning from data in psychology and machine learning” in *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society* (2022), pp. 335–348.
14. A. L. Beam, A. K. Manrai, M. Ghassemi, Challenges to the reproducibility of machine learning models in health care. *JAMA* **323**, 305–306 (2020).

15. M. Roberts, D. Driggs, M. Thorpe, J. Gilbey, M. Yeung, S. Ursprung, A. I. Aviles-Rivero, C. Etmann, C. McCague, L. Beer, J. R. Weir-McCall, Z. Teng, E. Gkrania-Klotsas, J. H. F. Rudd, E. Sala, C.-B. Schonlieb, Common pitfalls and recommendations for using “ machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans. *Nat. Mach. Intell.* **3**, 199–217 (2021).
16. X. Bouthillier, C. Laurent, P. Vincent, “Unreproducible research is reproducible” in *International Conference on Machine Learning* (PMLR, 2019), pp. 725–734.
17. M. B. McDermott, S. Wang, N. Marinsek, R. Ranganath, L. Foschini, M. Ghassemi, Reproducibility in machine learning for health research: Still a ways to go. *Sci. Transl. Med.* **13**, eabb1655 (2021).
18. O. E. Gundersen, S. Kjensmo, State of the art: Reproducibility in artificial Intelligence. *Proc. Conf. Artif. Intell.* **32** (2018).
19. G. Varoquaux, V. Cheplygina, Machine learning for medical imaging: Methodological failures and recommendations for the future. *NPJ Digit. Med.* **5**, 48 (2022).
20. L. Messeri, M. J. Crockett, Artificial intelligence and illusions of understanding in scientific research. *Nature* **627**, 49–58 (2024).
21. R. M. Schmidt, F. Schneider, P. Hennig, “Descending through a crowded valley - benchmarking deep learning optimizers” in *Proceedings of the 38th International Conference on Machine Learning* (PMLR, 2021), pp. 9367–9376.
22. X. Bouthillier, P. Delaunay, M. Bronzi, A. Trofimov, B. Nichyporuk, J. Szeto, N. Mohammadi Sepahvand, E. Raff, K. Madan, V. Voleti, S. E. Kahou, V. Michalski, T. Arbel, C. Pal, G. Varoquaux, P. Vincent, Accounting for variance in machine learning benchmarks, *Proc. Mach. Learn. Syst.* **3**, 747–769 (2021).
23. O. DeMasi, K. Kording B. Recht, Meaningless comparisons lead to false optimism in medical machine learning, *PLOS ONE* **12**, e0184604 (2017).
24. D. Sculley, J. Snoek, A. Wiltschko, A. Rahimi, *Winner’s curse? On pace, progress, and empirical rigor* (2018); <https://openreview.net/forum?id=rJWF0Fywf>.
25. D. M. Liu, M. J. Salganik, Successes and struggles with computational reproducibility: Lessons from the fragile families challenge. *Socius* **5**, 10.1177/2378023119849803 (2019).
26. S. Dittmer, M. Roberts, J. Gilbey, A. Biguri, AIX-COVNET Collaboration, J. Preller, J. H. F. Rudd, J. A. D. Aston, C.-B. Schonlieb, Navigating the development challenges in creating complex data systems. *Nat. Mach. Intell.* **5**, 681–686 (2023).

27. O. E. Gundersen, S. Shamsaliei, R. J. Isdahl, Do machine learning platforms provide out-of-the-box reproducibility?. *Future Gener. Comput. Syst.* **126**, 34–47 (2022).
28. J. M. Hofman, A. Sharma, D. J. Watts, Prediction and explanation in social systems, *Science* **355**, 486–488 (2017).
29. J. Banja, AI hype and radiology: A plea for realism and accuracy. *Radiol. Artif. Intell.* **2**, e190223 (2020).
30. V. E. Johnson, R. D. Payne, T. Wang, A. Asher, S. Mandal, On the reproducibility of psychological science. *J. Am. Stat. Assoc.* **112**, 1–10 (2017).
31. S. J. Bell, O. P. Kampman, Perspectives on machine learning from psychology’s reproducibility crisis. arXiv:2104.08878 [cs.LG] (18 April 2021).
32. J. Pineau, P. Vincent-Lamarre, K. Sinha, V. Lariviere, A. Beygelzimer, F. d’Alche-Buc, E. Fox, H. Larochelle, Improving reproducibility in machine learning research (a report from the NeurIPS 2019 reproducibility program). *J. Mach. Learn. Res.* **22**, 7459–7478 (2022).
33. M. Serra-Garcia, U. Gneezy, Nonreplicable publications are cited more than replicable ones. *Sci. Adv.* **7**, eabd1705 (2021).
34. M. J. Salganik, I. Lundberg, A. T. Kindel, C. E. Ahearn, K. al-Ghoneim, A. Almaatouq, D. M. Altschul, J. E. Brand, N. B. Carnegie, R. J. Compton, D. Datta, T. Davidson, A. Filippova, C. Gilroy, B. J. Goode, E. Jahani, R. Kashyap, A. Kirchner, S. McKay, A. C. Morgan, A. Pentland, K. Polimis, L. Raes, D. E. Rigobon, C. V. Roberts, D. M. Stanescu, Y. Suhara, A. Usmani, E. H. Wang, M. Adem, A. Alhajri, B. AlShebli, R. Amin, R. B. Amos, L. P. Argyle, L. Baer-Bositis, M. Büchi, B.R. Chung, W. Eggert, G. Faletto, Z. Fan, J. Freese, T. Gadgil, J. Gagné, Y. Gao, A. Halpern-Manners, S. P. Hashim, S. Hausen, G. He, K. Higuera, B. Hogan, I. M. Horwitz, L. M. Hummel, N. Jain, K. Jin, D. Jurgens, P. Kaminski, A. Karapetyan, E. H. Kim, B. Leizman, N. Liu, M. Möser, A. E. Mack, M. Mahajan, N. Mandell, H. Marahrens, D. Mercado-Garcia, V. Mocz, K. Mueller-Gastell, A. Musse, Q. Niu, W. Nowak, H. Omidvar, A. Or, K. Ouyang, K. M. Pinto, E. Porter, K. E. Porter, C. Qian, T. Rauf, A. Sargsyan, T. Schaffner, L. Schnabel, B. Schonfeld, B. Sender, J. D. Tang, E. Tsurkov, A. van Loon, O. Varol, X. Wang, Z. Wang, J. Wang, F. Wang, S. Weissman, K. Whitaker, M. K. Wolters, W. L. Woon, J. Wu, C. Wu, K. Yang, J. Yin, B. Zhao, C. Zhu, J. Brooks-Gunn, B. E. Engelhardt, M. Hardt, D. Knox, K. Levy, A. Narayanan, B. M. Stewart, D. J. Watts, S. McLanahan, Measuring the predictability of life outcomes with a scientific mass collaboration. *Proc. Natl. Acad. Sci. U.S.A.* **117**, 8398–8403 (2020).

35. J. Mongan, L. Moy, C. E. Kahn, Checklist for artificial intelligence in medical imaging (CLAIM): A guide for authors and reviewers. *Radiol. Artif. Intell.* **2**, e200029 (2020).
36. P. M. Bossuyt, J. B. Reitsma, D. E. Bruns, C. A. Gatsonis, P. P. Glasziou, L. Irwig, J. G. Lijmer, D. Moher, D. Rennie, H. C. W. de Vet, H. Y. Kressel, N. Rifai, R. M. Golub, D. G. Altman, L. Hooft, D. A. Korevaar, J. F. Cohen, STARD 2015: An updated list of essential items for reporting diagnostic accuracy studies. *BMJ Open* **351**, h5527 (2015).
37. R. G. White, A. J. Hakim, M. J. Salganik, M. W. Spiller, L. G. Johnston, L. Kerr, C. Kendall, A. Drake, D. Wilson, K. Orroth, M. Egger, W. Hladik, Strengthening the Reporting of Observational Studies in Epidemiology for respondent-driven sampling studies: STROBE-RDS statement. *J. Clin. Epidemiol.* **68**, 1463–1471 (2015).
38. E. G. Simmonds, K. P. Adjei, C. W. Andersen, Janne Cathrin Hetle Aspheim, C. Battistin, N. Bulso, H. Christensen, B. Cretois, R. Cubero, I. A. Davidovich, L. Dickel, B. Dunn, E. Dunn-Sigouin, K. Dyrstad, S. Einum, D. Giglio, H. Gjerlow, A. Godefroidt, R. Gonzalez-Gil, S. G. Cogno, F. Grosse, P. Halloran, M. F. Jensen, J. J. Kennedy, P. E. Langsaether, J. H. Laverick, D. Lederberger, C. Li, E. Mandeville, C. Mandeville, E. Moe, T. N. Schroder, D. Nunan, J. S. Parada, M. R. Simpson, E. S. Skarstein, C. Spensberger, R. Stevens, A. Subramanian, L. Svendsen, O. M. Theisen, C. Watret, R. B. OHara, How is model-related uncertainty quantified and reported in different disciplines? arXiv:2206.12179 [stat.AP] (24 June 2022).
39. S. Kapoor, A. Narayanan, Leakage and the reproducibility crisis in machine-learningbased science. *Patterns* **4**, 100804 (2023).
40. G. S. Collins, J. B. Reitsma, D. G. Altman, K. G. Moons, Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis (TRIPOD): The TRIPOD statement. *BMJ* **350**, g7594 2015.
41. A. C. Plint, D. Moher, A. Morrison, K. Schulz, D. G. Altman, C. Hill, I. Gaboury, Does the CONSORT checklist improve the quality of reports of randomised controlled trials? A systematic review. *Med. J. Aust.* **185**, 263–267 (2006).
42. S. Han, T. F. Olonisakin, J. P. Pribis, J. Zupetic, J. H. Yoon, K. M. Holleran, K. Jeong, N. Shaikh, D. M. Rubio, J. S. Lee, A checklist is associated with increased quality of reporting preclinical biomedical research: A systematic review. *PLOS ONE* **12**, e0183591 (2017).
43. Principles and guidelines for reporting preclinical research (2015); www.nih.gov/research-training/rigor-reproducibility/principles-guidelinesreporting-preclinical-research.

44. Reporting guidelines. The EQUATOR Network; www.equatornetwork.org/reporting-guidelines/.
45. A. Rogers, T. Baldwin, K. Leins, ‘just what do you think you’re doing, Dave?’ A checklist for responsible data use in NLP, in *Findings of the Association for Computational Linguistics: EMNLP 2021*, M.-F. Moens, X. Huang, L. Specia, S. Wen-tau Yih, Eds. (Punta Cana, Dominican Republic), (Association for Computational Linguistics, 2021). pp. 4821–4833.
46. O. E. Gundersen, Y. Gil, D. W. Aha, On reproducible AI: Towards reproducible research, open science, and digital scholarship in AI publications. *AI Mag.* **39**, 56–68 (2018).
47. D. Donoho, 50 years of data science. *J. Comput. Graph. Stat.* **26**, 745–766 (2017).
48. J. M. Hofman, D. J. Watts, S. Athey, F. Garip, T. L. Griffiths, J. Kleinberg, H. Margetts, S. Mullainathan, M. J. Salganik, S. Vazire, A. Vespignani, T. Yarkoni, Integrating explanation and prediction in computational social science. *Nature* **595**, 181–188 (2021).
49. E. Winsberg, *Science in the Age of Computer Simulation* (University of Chicago Press, 2010).
50. J. Pfeffer M. M. Malik, “Simulating the dynamics of socio-economic systems” in *Networked Governance: New Research Perspectives*, B. Hollstein, W. Matiaske, K.-U. Schnapp, Eds. (Springer International Publishing, 2017), pp.. 143–161.
51. R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, E. Brynjolfsson, S. Buch, D. Card, R. Castellon, N. S. Chatterji, A. S. Chen, K. A. Creel, J. Davis, D. Demszky, C. Donahue, M. Doumbouya, E. Durmus, S. Ermon, J. Etchemendy, K. Ethayarajh, L. Fei-Fei, C. Finn, T. Gale, L. Gillespie, K. Goel, N. D. Goodman, S. Grossman, N. Guha, T. Hashimoto, P. Henderson, J. Hewitt, D. E. Ho, J. Hong, K. Hsu, J. Huang, T. F. Icard, S. Jain, D. Jurafsky, P. Kalluri, S. Karamcheti, G. Keeling, F. Khani, O. Khattab, P. W. Koh, M. S. Krass, R. Krishna, R. Kuditipudi, A. Kumar, F. Ladhak, M. Lee, T. Lee, J. Leskovec, I. Levent, X. L. Li, X. Li, T. Ma, A. Malik, C. D. Manning, S. Mirchandani, E. Mitchell, Z. Munitykwa, S. Nair, A. Narayan, D. Narayanan, B. Newman, A. Nie, J. C. Niebles, H. Nilforoshan, J. F. Nyarko, G. Ogut, L. J. Orr, I. Papadimitriou, J. S. Park, C. Piech, E. Portelance, C. Potts, A. Raghunathan, R. Reich, H. Ren, F. Rong, Y. H. Roohani, C. Ruiz, J. Ryan, C. R’e, D. Sadigh, S. Sagawa, K. Santhanam, A. Shih, K. P. Srinivasan, A. Tamkin, R. Taori, A. W. Thomas, F. Tramer, R. E. Wang, W. Wang, B. Wu, J. Wu, ` Y. Wu, S. M. Xie, M. Yasunaga, J. You, M. A. Zaharia, M. Zhang, T. Zhang, X. Zhang, Y. Zhang, L. Zheng, K. Zhou, P. Liang, On the opportunities and risks of foundation models. arXiv:2108.07258 [cs.LG] (16 August 2021).

52. S. Kapoor, A. Narayanan, OpenAI's policies hinder reproducible research on language models (2022); www.aisnakeoil.com/p/openais-policies-hinderreproducible.
53. L. Chen, M. Zaharia, J. Zou, How is ChatGPT's behavior changing over time? arXiv:2307.09009 [cs.CL] (18 July 2023).
54. J. Carrasquilla, R. G. Melko, Machine learning phases of matter. *Nat. Phys.* **13**, 431–434 (2017).
55. I. Lundberg, R. Johnson, B. M. Stewart, What is your estimand? Defining the target quantity connects statistical evidence to theory. *Am. Sociol. Review* **86**, 532–565 (2021).
56. R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, F. A. Wichmann, Shortcut learning in deep neural networks. *Nat. Mach. Intell.* **2**, 665–673 (2020).
57. L. Wilkinson, Statistical methods in psychology journals: Guidelines and explanations, *Am. Psychol.* **54**, 594–604 (1999).
58. A. Casteel, N. Bridier, Describing populations and samples in doctoral student research, *Int. J. Dr. Stud.* **16**, 339–362 (2021).
59. L. Tooth, R. Ware, C. Bain, D. M. Purdie, A. Dobson, Quality of reporting of observational longitudinal research, *Am. J. Epidemiol.* **161**, 280–288 (2005).
60. D. J. Simons, Y. Shoda, D. S. Lindsay, Constraints on generality (COG): A proposed addition to all empirical papers. *Perspect. Psychol. Sci.* **12**, 1123–1128 (2017).
61. J. Grimmer, M. E. Roberts, B. Stewart, *Text as data: A New framework for Machine Learning and the Social Sciences*. (Princeton Univ. Press, 2022).
62. I. Lundberg, J. E. Brand, N. Jeon, Researcher reasoning meets computational capacity: Machine learning for social science. *Soc. Sci. Res.* **108**, 102807 (2022).
63. A. Neufeld, D. Witten, Discussion of Breiman's 'two cultures': From two cultures to one. *Obs. Stud.* **7**, 171–174 (2021).
64. G. Shmueli, Comment on Breiman's "Two Cultures" (2002): From two cultures to multicultural. *Obs. Stud.* **7**, 197–201 (2021).
65. M. Baiocchi, J. Rodu, Reasoning using data: Two old ways and one new. *Obs. Stud.* **7**, 3–12 (2021).
66. E. L. Ogburn I. Shpitser, Causal modelling: The two cultures. *Obs. Stud.* **7**, 179–183 (2021).
67. M. Molina, F. Garip, Machine learning for sociology. *Annu. Rev. Sociol.* **45**, 27–45 (2019).
68. H. H. Rashidi, N. Tran, S. Albahra, L. T. Dang, Machine learning in health care and laboratory medicine: General overview of supervised learning and auto-ML. *Int. J. Lab. Hematol.* **43**, 15–22 (2021).

69. A. L. Beam, I. S. Kohane, Big data and machine learning in health care. *JAMA* **319**, 1317 (2018).
70. L. G. McCoy, C. T. Brenna, S. S. Chen, K. Vold, S. Das, Believing in black boxes: machine learning for healthcare does not need explainability to be evidence-based. *J. Clin. Epidemiol.* **142**, 252–257 (2022).
71. A. L. Tarca, V. J. Carey, X.-W. Chen, R. Romero, S. Draghici, Machine learning and its applications to Biology. *PLoS Comput. Biol.* **3**, e116 (2007).
72. V. Stodden, Reproducing statistical results. *Annu. Rev. Stat. Its Appl.* **2**, 1–19 (2015).
73. T. Herndon, M. Ash, R. Pollin, Does high public debt consistently stifle economic growth? A critique of Reinhart and Rogoff. *Camb. J. Econ.* **38**, 257–279 (2014).
74. R. Herzog, P. A. M. Mediano, F. E. Rosas, R. Carhart-Harris, Y. S. Perl, E. Tagliazucchi, R. Cofre, Retraction note: A mechanistic model of the neural entropy increase elicited by psychedelic drugs. *Sci. Rep.* **12**, 15500 (2022).
75. M. R. Berenbaum, Retraction for Shu et al., Signing at the beginning makes ethics salient and decreases dishonest self-reports in comparison to signing at the end. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2115397118 (2021).
76. M. Neunhoeffer, S. Sternberg, How cross-validation can go wrong and what to do about it. *Polit. Anal.* **27**, 101–106 (2019).
77. J. M. Hofman, D. G. Goldstein, S. Sen, F. Poursabzi-Sangdeh, J. Allen, L. L. Dong, B. Fried, H. Gaur, A. Hoq, E. Mbazor, N. Moreira, C. Muso, E. Rapp, R. Terrero, Expanding the scope of reproducibility research through data analysis replications. *Organ. Behav. Hum. Decis. Process* **164**, 192–202 (2021).
78. G. Vandewiele, I. Dehaene, G. Kovacs, L. Sterckx, O. Janssens, F. Ongenae, F. De Backere, F. De Turck, K. Roelens, J. Decruyenaere, S. Van Hoecke, T. Demeester, Overly optimistic prediction results on imbalanced data: A case study of flaws and benefits when applying over-sampling. *Artif. Intell. Med.* **111**, 101987 (2021).
79. J. P. A. Ioannidis, D. B. Allison, C. A. Ball, I. Coulibaly, X. Cui, A. C. Culhane, M. Falchi, C. Furlanello, L. Game, G. Jurman, J. Mangion, T. Mehta, M. Nitzberg, G. P. Page, E. Petretto, V. van Noort, Repeatability of published microarray gene expression analyses. *Nat. Genet.* **41**, 149–155 (2009).
80. T. Verstynen, K. P. Kording, Overfitting to ‘predict’ suicidal ideation. *Nat. Hum. Behav.* **7**, 680–681 (2023).

81. B. Haibe-Kains, G. A. Adam, A. Hosny, F. Khodakarami, L. Waldron, B. Wang, C. McIntosh, A. Goldenberg, A. Kundaje, C. S. Greene, T. Broderick, M. M. Hoffman, J. T. Leek, K. Korthauer, W. Huber, A. Brazma, J. Pineau, R. Tibshirani, T. Hastie, J. P. A. Ioannidis, J. Quackenbush, H. J. W. L. Aerts, Transparency and reproducibility in artificial intelligence. *Nature* **586**, E14–E16 (2020).
82. V. Stodden, J. Seiler, Z. Ma, An empirical analysis of journal policy effectiveness for computational reproducibility. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 2584–2589 (2018).
83. M. Gabelica, R. Bojicic, L. Puljak, Many researchers were not compliant with their published data sharing statement: A mixed-methods study. *J. Clin. Epidemiol.* **150**, 33–41 (2022).
84. N. A. Vasilevsky, J. Minnier, M. A. Haendel, R. E. Champieux, Reproducible and reusable research: Are journal data sharing policies meeting the mark?. *PeerJ* **5**, e3208 (2017).
85. P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, D. Meger, Deep reinforcement learning that matters. *Proc. AAAI Conf. Artif. Intell.* **32**, 3207–3214 (2018).
86. K. Musgrave, S. Belongie, S.-N. Lim, “A metric learning reality check” in *Computer Vision—ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, J.-M. Frahm, Eds. (Springer International Publishing, Cham, 2020), pp. 681–699.
87. AAAI, AAAI reproducibility checklist; <https://aaai.org/conference/aaai/aaai23/reproducibility-checklist/>.
88. NeurIPS, NeurIPS 2023 paper guidelines; <https://neurips.cc/public/guides/PaperChecklist>.
89. ICML, ICML 2023 paper guidelines; <https://icml.cc/Conferences/2023/PaperGuidelines>.
90. The Journal of Politics, Guidelines for data replication;
www.journals.uchicago.edu/journals/jop/data-replication.
91. Science, Science journals: Editorial policies; www.science.org/content/page/sciencejournals-editorial-policies.
92. The Journal of Politics, Guidelines for data replication;
www.journals.uchicago.edu/journals/jop/data-replication.
93. B. A. Nosek, G. Alter, G. C. Banks, D. Borsboom, S. D. Bowman, S. J. Breckler, S. Buck, C. D. Chambers, G. Chin, G. Christensen, M. Contestabile, A. Dafoe, E. Eich, J. Freese, R. Glennerster, D. Goroff, D. P. Green, B. Hesse, M. Humphreys, J. Ishiyama, D. Karlan, A. Kraut, A. Lupia, P. Mabry, T. Madon, N. Malhotra, E. Mayo-Wilson, M. McNutt, E. Miguel, E. L. Paluck, U. Simonsohn, C.

- Soderberg, B. A. Spellman, J. Turitto, G. VandenBos, S. Vazire, E. J. Wagenmakers, R. Wilson, T. Yarkoni, Promoting an open research culture. *Science* **348**, 1422–1425 (2015).
94. Koren, Miklos, Connolly, Marie, Lull, Joan, Vilhuber, Lars, Data and code availability standard (2022); <https://zenodo.org/record/7436134>.
95. Reviewing computational methods. *Nat. Methods* **12**, 1099–1099 (2015).
96. K. Peng, A. Mathur, A. Narayanan, Mitigating dataset harms requires stewardship: Lessons from 1000 papers. *Proc. Neural Inf. Process. Syst. Track on Datasets Benchmarks*, vol. 1, 2021.
97. U.S. Geological Survey, Data dictionaries; www.usgs.gov/datamanagement/data-dictionaries.
98. T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, K. Crawford, Datasheets for datasets, *Commun. ACM* **64**, 86–92 (2021).
99. J. Walonoski, M. Kramer, J. Nichols, A. Quina, C. Moesel, D. Hall, C. Duffett, K. Dube, T. Gallagher, S. McLachlan, Synthea: An approach, method, and software mechanism for generating synthetic patients and the synthetic electronic health care record. *J. Am. Med. Inform. Assoc.* **25**, 230–238 (2018).
100. J. Jordon, L. Szpruch, F. Houssiau, M. Bottarelli, G. Cherubin, C. Maple, S. N. Cohen, A. Weller, Synthetic Data—what, why and how? arXiv:2205.03257 [cs.LG] (6 May 2022).
101. Z. Obermeyer, B. Powers, C. Vogeli, S. Mullainathan, Dissecting racial bias in an algorithm used to manage the health of populations. *Science* **366**, 447–453 (2019).
102. B. Nowok, G. M. Raab, C. Dibben, Synthpop: Bespoke creation of synthetic data in R. *J. Stat. Softw.* **74**, 1–26, (2016).
103. J. Chen, D. Chun, M. Patel, E. Chiang, J. James, The validity of synthetic clinical data: a validation study of a leading synthetic data generator (Synthea) using clinical quality measures, *BMC Med. Inform. Decis. Mak.* **19**, 44 (2019).
104. S. Hao, W. Han, T. Jiang, Y. Li, H. Wu, C. Zhong, Z. Zhou, H. Tang, Synthetic data in AI: Challenges, applications, and ethical implications. arXiv:2401.01629v1 [cs.LG] (3 January 2024).
105. G. K. Sandve, A. Nekrutenko, J. Taylor, E. Hovig, Ten simple rules for reproducible computational Research. *PLoS Comput. Biol.* **9**, e1003285 (2013).
106. T. H. Vines, A. Y. Albert, R. L. Andrew, F. Debarre, D. G. Bock, M. T. Franklin, K. J. Gilbert, J.-S. Moore, S. Renaut, D. J. Rennison, The availability of research data declines rapidly with article age, *Curr. Biol.* **24**, 94–97 (2014).

107. E. Gibney, R. Van Noorden, Scientists losing data at a rapid rate, *Nature*, (2013).
<https://doi.org/10.1038/nature.2013.14416>.
108. V. Stodden, S. Miguez, Best practices for computational science: Software infrastructure and environments for reproducible and extensible research, *J. Open Research Softw.* **2**, e21 (2014).
109. L. Vilhuber, M. Connolly, M. Koren, J. Llull, P. Morrow, A template README for social science replication packages (2020); <https://zenodo.org/record/4319999>.
110. M. Singers, Awesome README; <https://github.com/matiassingers/awesome-readme>.
111. Harbert, Bash scripting (2018); <https://rsh249.github.io/bioinformatics/bash-script.html>.
112. A. Clyburne-Sherin, X. Fei, S. A. Green, Computational reproducibility via containers in psychology. *Meta Psychol.* **3** (2019).
113. C. L. A. Navarro, J. A. A. Damen, T. Takada, S. W. J. Nijman, P. Dhiman, J. Ma, G. S. Collins, R. Bajpai, R. D. Riley, K. G. M. Moons, L. Hooft, Completeness of reporting of clinical prediction models developed using supervised machine learning: A systematic review. *BMC Med. Res. Methodol.* **22**, 12 (2022).
114. M. Yusuf, I. Atal, J. Li, P. Smith, P. Ravaud, M. Fergie, M. Callaghan, J. Selfe, Reporting quality of studies using machine learning models for medical diagnosis: A systematic review. *BMJ Open* **10**, e034568 (2020).
115. Y. Kim, J. Huang, S. Emery, Garbage in, garbage out: Data collection, quality assessment and reporting standards for social media data use in health research, infodemiology and digital disease detection. *J. Med. Internet Res.* **18**, e41 (2016).
116. R. S. Geiger, K. Yu, Y. Yang, M. Dai, J. Qiu, R. Tang, J. Huang, “Garbage in, garbage out? do machine learning application papers in social computing report where humanlabeled training data comes from?” in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, FAT* '20, (Association for Computing Machinery, 2020).
117. F. Porzsolt, F. Wiedemann, S. I. Becker, C. J. Rhoads, Inclusion and exclusion criteria and the problem of describing homogeneity of study populations in clinical trials. *BMJ Evid. Based Med.* **24**, 92–94 (2019).
118. M. Salganik, *Bit by Bit: Social Research in the Digital Age* (Princeton Univ. Press, 2019).
119. S. Barocas, M. Hardt, A. Narayanan, *Fairness and Machine Learning: Limitations and Opportunities* (2019); www.fairmlbook.org.

120. A. Z. Jacobs, H. Wallach, “Measurement and fairness” in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21 (Association for Computing Machinery, New York, NY, USA, 2021), pp. 375–385.
121. M. Crede, P. D. Harms, Three cheers for descriptive statistics—and five more reasons why they matter. *Ind. Organ. Psychol.* **14**, 486–488 (2021).
122. J. Larson-Hall, L. Plonsky, Reporting and interpreting quantitative research findings: What gets reported and recommendations for the field. *Lang. Learn.* **65**, 127–159 (2015).
123. V. C. Bradley, S. Kuriwaki, M. Isakov, D. Sejdinovic, X.-L. Meng, S. Flaxman, Unrepresentative big surveys significantly overestimated US vaccine uptake. *Nature* **600**, 695–700 (2021).
124. L. Plonsky, Study quality in SLA. *Stud. Second Lang. Acquis.* **35**, 655–687 (2013).
125. P. McKnight, K. McKnight, S. Sidani, A. J. Figueredo, *Missing Data: A Gentle Introduction* (Guilford Publications, 2007).
126. J. L. Peugh, C. K. Enders, Missing data in educational research: A review of reporting practices and suggestions for improvement, *Rev. Educ. Res.* **74**, 525–556 (2004).
127. C. Mack, Z. Su, D. Westreich, *Managing Missing Data in Patient Registries: Addendum to Registries for Evaluating Patient Outcomes: A User’s Guide (Third edition)* (2018); www.ncbi.nlm.nih.gov/books/NBK493611/.
128. S. Nijman, A. Leeuwenberg, I. Beekers, I. Verkouter, J. Jacobs, M. Bots, F. Asselbergs, K. Moons, T. Debray, Missing data is poorly handled and reported in prediction model studies using machine learning: A literature review. *J. Clin. Epidemiol.* **142**, 218–229 (2022).
129. T. D. Little, T. D. Jorgensen, K. M. Lang, E. W. G. Moore, On the joys of missing data. *J. Pediatr. Psychol.* **39**, 151–162 (2014).
130. J. S. Nicholson, P. R. Deboeck, W. Howard, Attrition in developmental psychology. *Int. J. Behav. Dev.* **41**, 143–153 (2017).
131. W. R. Sterner, What is missing in counseling research? reporting missing data. *J. Couns. Dev.* **89**, 56–62 (2011).
132. J. A. Hussain, M. Bland, D. Langan, M. J. Johnson, D. C. Currow, I. R. White, Quality of missing data reporting and handling in palliative care trials demonstrates that further development of the CONSORT statement is required: a systematic review. *J. Clin. Epidemiol.* **88**, 81–91 (2017).
133. X. Chu, I. F. Ilyas, S. Krishnan, J. Wang, “Data cleaning: Overview and emerging challenges” in *Proceedings of the 2016 International Conference on Management of Data* (2016); pp. 2201–2206.

134. E. M. Buchanan, J. E. Scofield, Methods to detect low quality data and its implication for psychological research. *Behav. Res. Methods* **50**, 2586–2596 (2018).
135. T. Shadbahr, M. Roberts, J. Stanczuk, J. Gilbey, P. Teare, S. Dittmer, M. Thorpe, R. V. Torne, E. Sala, P. Lio, M. Patel, AIX-COCNET Collaboration; J. H. F. Rudd, T. Mirtti, A. Rannikko, J. A. D. Aston, J. Tang, C.-B. Schönlieb, Classification of datasets with imputed missing values: Does imputation quality matter? arXiv: 2206.08478 [cs.LG] (16 June 2022).
136. E. Gryska, I. Bjorkman-Burtscher, A. S. Jakola, T. Dunas, J. Schneiderman, R. A. Heckemann, Deep learning for automatic brain tumour segmentation on mri: Evaluation of recommended reporting criteria via a reproduction and replication study, *BMJ Open* **12**, e059000, (2022).
137. E. Raff, “A step toward quantifying independently reproducible machine learning research” in *Advances in Neural Information Processing Systems*, vol. 32 (2019), pp. 5485–5495.
138. M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, T. Gebru, “Model cards for model reporting” in *Proceedings of the Conference on Fairness, Accountability, and Transparency* (2019), pp. 220–229.
139. J. Kleinberg, A. Liang, S. Mullainathan, “The theory is predictive, but is it complete? An application to human perception of randomness” in *Proceedings of the 2017 ACM Conference on Economics and Computation* (2017), pp. 125–126.
140. R. Roscher, B. Bohn, M. F. Duarte, J. Garcke, Explainable machine learning for scientific insights and discoveries. *IEEE Access* **8**, pp. 42200–42216 (2020).
141. C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* **1**, 206–215 (2019).
142. S. Raschka, Model evaluation, model selection, and algorithm selection in machine learning. arXiv:1811.12808v3 [cs.LG] (11 November 2020).
143. G. C. Cawley, N. L. C. Talbot, On over-fitting in model selection and subsequent selection bias in performance evaluation. *J. Mach. Learn. Res.* **11**, 2079–2107 (2010).
144. S. Wager, Cross-validation, risk estimation, and model selection: Comment on a paper by Rosset and Tibshirani. *J. Am. Stat. Assoc.* **115**, 157–160 (2020).
145. C. Marx, F. Calmon, B. Ustun, “Predictive multiplicity in classification” in *Proceedings of the 37th International Conference on Machine Learning* (PMLR, 2020), pp. 6765–6774.
146. J. Watson-Daniels, D. C. Parkes, B. Ustun, Predictive multiplicity in probabilistic classification. *Proc AAAI Conf. Artif. Intell.* **37**, 10306–10314 (2023).

147. E. Black, M. Raghavan, S. Barocas, Model multiplicity: Opportunities, concerns, and solutions, in 2022 ACM Conference on Fairness, Accountability, and Transparency, Seoul, Republic of Korea, 2022 June 21 to 24, 2022.
148. J. Dodge, S. Gururangan, D. Card, R. Schwartz, N. A. Smith, Show your work: Improved reporting of experimental results. arXiv:1909.03004 [cs.LG] (6 September 2019).
149. R. Islam, P. Henderson, M. Gomrokchi, D. Precup, Reproducibility of benchmarked deep reinforcement learning tasks for continuous control. arXiv:1708.04133 [cs.LG] (10 August 2017).
150. A. F. Cooper, Y. Lu, J. Forde, C. M. De Sa, “Hyperparameter optimization is deceiving us, and how to stop it” in *Advances in Neural Information Processing Systems* (Curran Associates, Inc. 2021), vol. 34, pp. 3081–3095.
151. P. T. Sivaprasad, F. Mai, T. Vogels, M. Jaggi, Fe. Fleuret, Optimizer benchmarking needs to account for hyperparameter tuning, in *Proceedings of the 37th International Conference on Machine Learning* (PMLR, 2020), pp. 9036–9045.
152. G. E. Dahl, F. Schneider, Z. Nado, N. Agarwal, C. S. Sastry, P. Hennig, S. Medapati, R. Eschenhagen, P. Kasimbeg, D. Suo, J. Bae, J. Gilmer, A. L. Peirson, B. Khan, R. Anil, M. Rabbat, S. Krishnan, D. Snider, E. Amid, K. Chen, C. J. Maddison, R. Vasudev, M. Badura, A. Garg, P. Mattson, Benchmarking neural network training algorithms. arXiv:2306.07179 [cs.LG] (12 June 2023).
153. P. Probst, A.-L. Boulesteix, B. Bischl, Tunability: Importance of hyperparameters of machine learning algorithms. *J. Mach. Learn. Res.* **20**, 1934–1965 (2019).
154. J. Lin, The neural hype and comparisons against weak baselines. *ACM SIGIR Forum* **52**, 40–51 (2019).
155. M. A. Lones, How to avoid machine learning pitfalls: A guide for academic researchers. arXiv:2108.02497v3 [cs.LG] (9 February 2023).
156. S. Kaufman, S. Rosset, C. Perlich, O. Stitelman, Leakage in data mining: Formulation, detection, and avoidance. *ACM Trans. Knowl. Discov. Data* **6**, 1–21 (2012).
157. C. Ross, Epic’s sepsis algorithm is going off the rails in the real world. The use of these variables may explain why, 2021; www.statnews.com/2021/09/27/epic-sepsisalgorithm-antibiotics-model/.
158. M. Kuhn, K. Johnson, *Applied Predictive Modeling* (Springer-Verlag, 2013).
159. R. A. Poldrack, G. Huckins, G. Varoquaux, Establishment of best practices for evidence for prediction: A review. *JAMA Psychiatry* **77**, 534–540 (2020).

160. M. U. Oner, Y.-C. Cheng, H. K. Lee, W.-K. Sung, *Training machine learning models on patient level data segregation is crucial in practical clinical applications* (2020); www.medrxiv.org/content/10.1101/2020.04.23.20076406v1.
161. M. M. Malik, A hierarchy of limitations in machine learning. arXiv:2002.05193 [cs.CY] (12 February 2020).
162. M. Lachanski, S. Pav, Shy of the character limit: Twitter mood predicts the stock market revisited. *Econ J. Watch* **14**, 302–345 (2017).
163. C. Bergmeir, J. M. Benitez, On the use of cross-validation for time series predictor evaluation. *Inform. Sci.* **191**, 192–213 (2012).
164. N. Y. Hammerla, T. Plotz, “Let’s (not) stick together: Pairwise similarity biases cross-validation in activity recognition” in *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, Osaka, Japan (ACM, 2015), pp. 1041–1051.
165. D. R. Roberts, V. Bahn, S. Ciuti, M. S. Boyce, J. Elith, G. Guillera-Aroita, S. Hauenstein, J. J. Lahoz-Monfort, B. Schroder, W. Thuiller, D. I. Warton, B. A. Wintle, F. Hartig, C. F. Dormann, Cross-validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure. *Ecography* **40**, 913–929 (2017).
166. A. Chiavegatto Filho, A. F. D. M. Batista, H. G. Dos Santos, Data leakage in health outcomes prediction with machine learning. Comment on Prediction of incident hypertension within the next year: Prospective study using statewide electronic health records and machine learning. *J. Med. Internet Res.* **23**, e10969 (2021).
167. C. Ye, T. Fu, S. Hao, Y. Zhang, O. Wang, B. Jin, M. Xia, M. Liu, X. Zhou, Q. Wu, Y. Guo, C. Zhu, Y.-M. Li, D. S. Culver, S. T. Alfreds, F. Stearns, K. G. Sylvester, E. Widen, D. McElhinney, X. Ling, Prediction of incident hypertension within the next year: Prospective study using statewide electronic health records and machine learning. *J. Med. Internet Res.* **20**, e22 (2018).
168. J. Z. Forde, A. F. Cooper, K. Kwegyir-Aggrey, C. De Sa, M. Littman, Model selection’s disparate impact in real-world deep learning applications. arXiv:2104.00606 [cs.LG] (1 April 2021).
169. U. Bhatt, J. Antoran, Y. Zhang, Q. V. Liao, P. Sattigeri, R. Fogliato, G. Melanc, on, R. Krishnan, J. Stanley, O. Tickoo, L. Nachman, R. Chunara, M. Srikumar, A. Weller, A. Xiang, “Uncertainty as a Form of Transparency: Measuring, Communicating, and Using Uncertainty” in *Proceedings of the*

- 2021 AAAI/ACM Conference on AI, Ethics, and Society, AIES '21 (Association for Computing Machinery, New York, NY, USA, 2021), pp. 401–413.
170. A. F. Cooper, K. Lee, M. Z. Choksi, S. Barocas, C. De Sa, J. Grimmelman, J. Kleinberg, S. Sen, B. Zhang, Arbitrariness and Social Prediction: The Confounding Role of Variance in Fair Classification. *Proc. AAAI Conf. Artif. Intell.* **38**, 22004–22012 (2024).
171. S. Qian, V. H. Pham, T. Lutellier, Z. Hu, J. Kim, L. Tan, Y. Yu, J. Chen, S. Shah, “Are my deep learning systems fair? An empirical study of fixed-seed training, in *Advances in Neural Information Processing Systems* (Curran Associates Inc. 2021), vol. 34, (pp. 30211–30227).
172. C. Young, Model Uncertainty and the Crisis in Science, *Socius* **4**, 2378023117737206 (2018).
173. A. N. Angelopoulos, S. Bates, C. Fannjiang, M. I. Jordan, T. Zrnic, Prediction-powered inference. *Science*. **382**, 669–674 (2023).
174. M. C. Monard, G. Batista, Learning with skewed class distributions. *Adv. Logic, Artif. Intell. Robot.* **85**, 173–180 (2002).
175. J. M. Lobo, A. Jimenez-Valverde, R. Real, Auc: A misleading measure of the performance of predictive distribution models, *Glob. Ecol. Biogeogr.* **17**, 145–151 (2008).
176. A. Bhowmick, S. M. Hazarika, E-mail spam filtering: A review of techniques and trends. *Adv. Electron. Commun. Comput.* **443**, 583–590 2018.
177. X.-H. Zhou, D. K. McClish, N. A. Obuchowski, *Statistical Methods in Diagnostic Medicine* (John Wiley & Sons, ed. 2, 2011).
178. V. Amrhein, S. Greenland, B. McShane, Scientists rise up against statistical significance. *Nature* **567**, 305–307 (2019).
179. W. R. Shadish, T. D. Cook, D. T. Campbell, *Experimental and Quasi-Experimental Designs for Generalized Causal Inference* (Houghton Mifflin, ed. 2, 2001).
180. N. Egami, E. Hartman, Elements of external validity: Framework, design, and analysis. *Am. Polit. Sci. Rev.* **117**, 1070–1088 (2022).
181. I. D. Raji, I. E. Kumar, A. Horowitz, A. Selbst, “The fallacy of AI functionality” in *2022 ACM Conference on Fairness, Accountability, and Transparency* (ACM, Seoul Republic of Korea, 2022), p. 959–972.
182. T. Liao, R. Taori, I. D. Raji, L. Schmidt, “Are we learning yet? a meta review of evaluation failures across machine learning” in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track* (2021).

183. A. Mathur, A. Wang, C. Schwemmer, M. Hamin, B. M. Stewart, A. Narayanan, Manipulative tactics are the norm in political emails: Evidence from 300k emails from the 2020 US election cycle. *Big Data Soc.* **10**, 205395172211453 (2023).
184. J. Brownlee, Difference between algorithm and model in machine learning (2020); <https://machinelearningmastery.com/difference-between-algorithm-and-model-in-machine-learning/>.
185. A. F. Dugas, M. Jalalpour, Y. Gel, S. Levin, F. Torcaso, T. Igusa, R. E. Rothman, Influenza forecasting with Google Flu trends, *PLOS ONE* **8**, e56176 (2013).
186. O. Wiles, S. Gowal, F. Stimberg, S. Alvisè-Rebuffi, I. Ktena, K. Dvijotham, T. Cemgil, A fine-grained analysis on distribution shift. arXiv:2110.11328 [cs.LG] (21 October 2021).
187. P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga, R. L. Phillips, I. Gao, T. Lee, E. David, I. Stavness, W. Guo, B. A. Earnshaw, I. S. Haque, S. Beery, J. Leskovec, A. Kundaje, E. Pierson, S. Levine, C. Finn, P. Liang, “Wilds: A benchmark of in-the-wild distribution shifts” in *International Conference on Machine Learning* (PMLR, 2021), pp. 5637–5664.
188. K. Bansak, J. Ferwerda, J. Hainmueller, A. Dillon, D. Hangartner, D. Lawrence, J. Weinstein, Improving refugee integration through data-driven algorithmic assignment. *Science* **359**, 325–329 (2018).
189. S. Bozkurt, E. M. Cahan, M. G. Seneviratne, R. Sun, J. A. Lossio-Ventura, J. P. A. Ioannidis, T. Hernandez-Boussard, Reporting of demographic data and representativeness in machine learning models using electronic health records, *J. Am. Med. Inform. Assoc.* **27**, 1878–1884 (2020).
190. C. L. A. Navarro, J. A. Damen, T. Takada, S. W. Nijman, P. Dhiman, J. Ma, G. S. Collins, R. Bajpai, R. D. Riley, K. G. Moons, L. Hooft, Systematic review finds spin practices and poor reporting standards in studies on machine learning-based prediction models. *J. Clin. Epidemiol.* **158**, 99–110 (2023).
191. S. G. Finlayson, A. Subbaswamy, K. Singh, J. Bowers, A. Kupke, J. Zittrain, I. S. Kohane, S. Saria, The clinician and dataset shift in artificial intelligence. *New Engl. J. Med.* **385**, 283–286 (2021).
192. M. Macleod, A. M. Collings, C. Graf, V. Kiermer, D. Mellor, S. Swaminathan, D. Sweet, V. Vinson, The MDAR (Materials Design Analysis Reporting) framework for transparent reporting in the life sciences. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2103238118 (2021).
193. Unofficial guidance on various topics by social science data editors; <https://socialscience-data-editors.github.io/guidance/>.

194. Requirements file format–pip documentation v23.0.1;
<https://pip.pypa.io/en/stable/reference/requirements-file-format/>.
195. Nature research code and software submission checklist; www.nature.com/documents/nr-software-policy.pdf, 2017.
196. T. Comi, Using CodeOcean for sharing reproducible research;
<https://rse.princeton.edu/2021/03/using-codeocean-for-sharing-reproducible-research/>.
197. K. S. Chmielinski, S. Newman, M. Taylor, J. Joseph, K. Thomas, J. Yurkofsky, Y. C. Qiu, The Dataset Nutrition Label (2nd Gen): Leveraging context to mitigate harms in artificial intelligence. arXiv:2201.03954 [cs.LG] (10 January 2022).
198. *Brain imaging data structure* (2023); <https://bids.neuroimaging.io/index>.
199. H. Taherdoost, *Sampling methods in research methodology*; How to choose a sampling technique for research, 2016.
200. J. P. Vandembroucke, E. Von Elm, D. G. Altman, P. C. Gøtzsche, C. D. Mulrow, S. J. Pocock, C. Poole, J. J. Schlesselman, M. Egger, STROBE Initiative, Strengthening the Reporting of Observational Studies in Epidemiology (STROBE): Explanation and elaboration, *Int. J. Surg.* **4**, e297 (2007).
201. 3.2. Tuning the hyper-parameters of an estimator; <https://scikitlearn/stable/modules/gridsearch.html>.
202. A. Vehtari, Cross-validation FAQ; <https://avehtari.github.io/modelselection/CVFAQ.html>.