

Supporting Information

Table S1 Selected R tools that can help with missing data problems and their potential application for use in biodiversity research.

R packages	Applications	Useful functions
<i>Exploring missing data</i>		
naniar	visualizing/exploring the missing data pattern	<i>mcar_test</i> – Little's missing completely at random (MCAR) test <i>vis_miss</i> – plot the missing data for all variables
occAssess	measure of the potential for bias in taxonomic, temporal, spatial, and environmental dimensions	<i>assessEnvBias</i> - assess whether data are sampled from a representative portion of environmental space in the spatial domain of interest <i>assessSpatialBias</i> – assess whether data resemble a random distribution in the geographic space of interest for inference <i>assessSpatialCov</i> – assess whether a representative portion of the spatial domain of interest has been sampled and whether the same portion of geographic space has been sampled over time
sambias	a Bayesian approach to estimate how sampling rates vary as a function of proximity to one or multiple bias factors	<i>calculate_bias</i> - calculating the bias effect of sampling bias due to geographic structures, such as the vicinity to cities, airports, rivers and roads
<i>Subsampling</i>		
base	Base R functions	<i>sample</i> - sample data with predefined inclusion probabilities specified with the prob argument
sampling	draw random samples using different sampling schemes	<i>balancedcluster</i> – selects a balanced cluster sample according to defined auxiliary variables <i>strata</i> - stratified sampling with unequal probabilities.

spatialEco	spatial data manipulation and modelling	<i>stratified.random</i> - creates a stratified random sample of an sp class object <i>stratified.distance</i> - draws a minimum, and optional maximum constrained, distance sub-sampling
spThin	Spatial thinning of species occurrence records	<i>thin</i> - returns a dataset with the maximum number of records for a given thinning distance
terra	spatial data manipulation and processing	<i>spatSample</i> – sample a SpatRaster, SpatVector or SpatExtent objcy
<i>Imputation</i>		
agTrend	modelling regional trends with missing data	<i>mcmc.aggregate</i> - a zero-inflated, nonparameteric model with a definable observation model, augmenting missing values before calculating regional abundances
INLA/ inlabru	fitting Bayesian models, especially useful for spatial models via its spatial mesh	<i>inla/bru</i> - fit a Bayesian model using Integrated Nested Laplace approximation <i>predict</i> – draw predictions from the fitted model, where the prediction data frame can be a SpatialPointsDataFrame object
LORI	imputation of missing count data	<i>lori</i> – impute missing count data using a large covariate set, including interactions, with a LASSO penalty
mice	multiple imputation by chained equations	<i>mice</i> – multiple imputation method that will generate plausible values for any missing data – in the response and in any covariates
Rjags JAGS nimble	fitting Bayesian models allowing for missing values in the response	<i>Jags/runMCMC</i> - fitting Bayesian models allowing for imputation of missing values in the response during model fitting (options available for missing values in covariates too)
rtrim	functions to calculate annual indices and trends of abundances	fit a GLM imputing missing values based on mean site and year effects, with optional covariates
<i>Weighting</i>		

survey svyr	range of functions for analysis of data from complex surveys, including fitting models with weights	<i>Svyglm</i> – generalized linear models with survey weights <i>postStratify</i> – function for post-stratification to match the joint distribution of the variables of the population
svrep	Analysis of replicate/boosted survey weights	<i>svyby_repwts</i> – compare estimates from different sets of weights
twang	functions to estimate propensity scores and weights	<i>ps</i> - gradient boosted trees to predict non-response from covariates <i>bal.table</i> – compare covariate values between sample and population

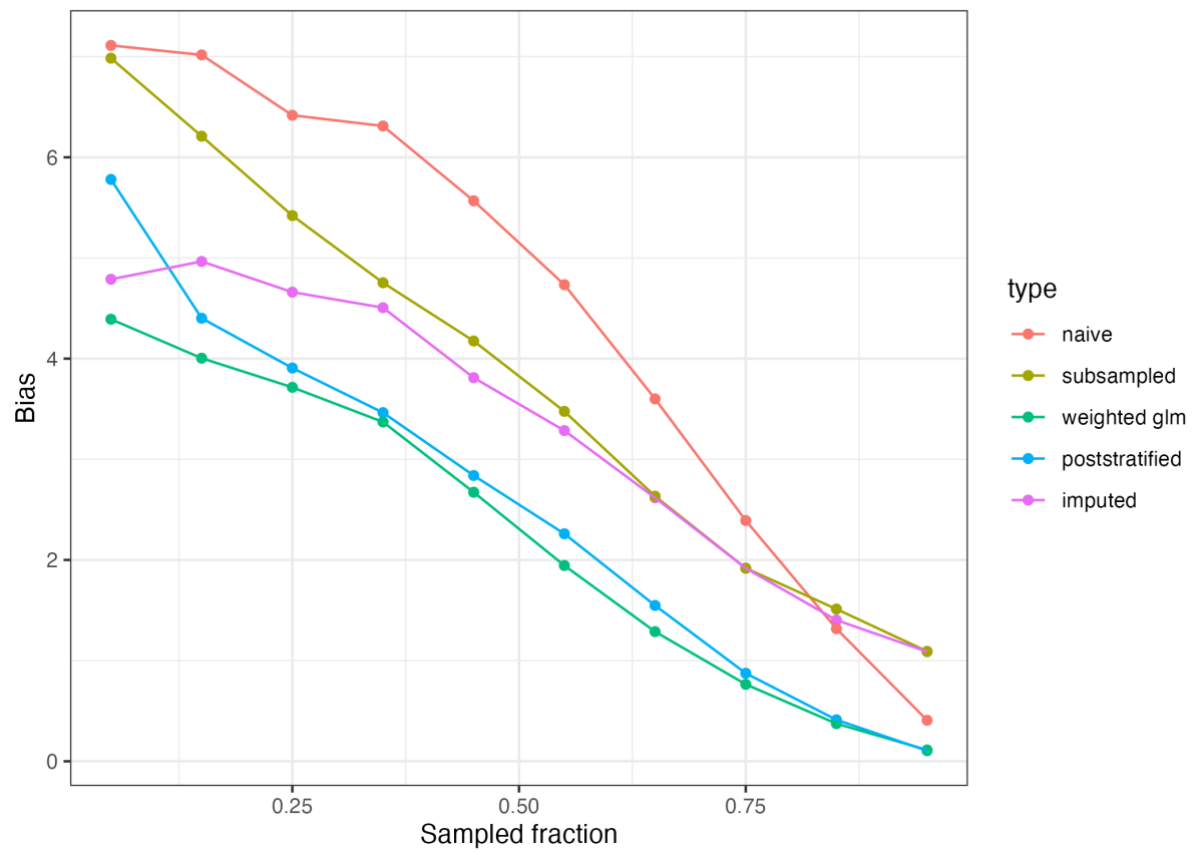


Fig. S1 The ability of missing data solutions to adjust for bias in biodiversity data.

We assumed a landscape of 400 cells and that a covariate affected both species abundance and the likelihood of a cell being sampled. We vary the fraction of the cells that were sampled. In contrast to Fig 5A (main text), we assumed that the species abundance was affected by an additional covariate that did not affect sampling; this variable was not included in any of the analysis. The models to estimate the parameter of interest (mean abundance) were: naive (no correction); subsampled (cells were subsampled along the covariate gradient to reduce the sampling bias), weighted (two methods: weighted glm using the svyglm function, and weighted by poststratification, using postStratify, both in the survey package) and imputed (using JAGS to impute NAs in the response). Points show the mean bias (difference between model prediction and truth) across 100 independent runs.