

1 **Reproducible, flexible and high throughput data extraction from primary**
2 **literature: The metaDigitise R package**

3 Joel L. Pick^{1,*}, Shinichi Nakagawa¹, Daniel W.A. Noble¹

4 ¹ Ecology and Evolution Research Centre, School of Biological, Earth and
5 Environmental Sciences, University of New South Wales, Kensington, NSW 2052,
6 Sydney, AUSTRALIA

7 *Corresponding Author: joel.l.pick@gmail.com

8 Abstract

- 9 1. Research synthesis, such as meta-analysis requires the extraction of effect sizes
10 from primary literature. Such effect sizes are calculated from summary statistics.
11 However, exact values of such statistics are commonly hidden in figures.
- 12 2. Extracting summary statistics from figures can be a slow process that is not easily
13 reproducible. Additionally, current software lacks an ability to incorporate
14 important meta-data (e.g., sample sizes, treatment / variable names) about
15 experiments and is not integrated with other software to streamline analysis
16 pipelines.
- 17 3. Here we present the R package **metaDigitise** which extracts descriptive statistics
18 such as means, standard deviations and correlations from the four plot types: 1)
19 mean/error plots (e.g. bar graphs with standard errors), 2) box plots, 3) scatter
20 plots and 4) histograms. **metaDigitise** is user-friendly and easy to learn as it
21 interactively guides the user through the data extraction process. Notably, it
22 enables large-scale extraction by automatically loading image files, letting the user
23 stop processing, edit and add to the resulting data frame at any point.
- 24 4. Digitised data can be easily re-plotted and checked, facilitating reproducible data
25 extraction from plots with little inter-observer bias. We hope that by making the
26 process of figure extraction more flexible and easy to conduct it will improve the
27 transparency and quality of meta-analyses in the future.

28 **Keywords:** meta-analysis, comparative analysis, data extraction, R, reproducibility,
29 figures, images, summary statistics

1 Introduction

In many different contexts, researchers make use of data presented in primary literature. In the ecology and evolution, this most notably includes comparative and meta-analyses. These techniques rely on descriptive statistics (e.g. means, standard deviations (SD), sample sizes, correlation coefficients) extracted from primary literature. As well as being presented in the text or tables of research papers, descriptive statistics are frequently presented in figures and so need to be manually extracted using digitising programs.

Although there are several tools that extract data from figures (e.g. **DataThief** (Tummers, 2006), **GraphClick** (Arizona-Software, 2008), **WebPlotDigitizer** (Rohatgi, 2017)), these tools do not cater to needs of meta-analysis for four main reasons (here we focus on meta-analysis, although many of the same points relate to comparative analysis). First, although meta-analysis is an important tool in consolidating the data from multiple studies, many of the processes involved in data extraction are opaque and difficult to reproduce, making extending or replicating studies problematic. Having a tool that facilitates reproducibility in meta-analyses will increase transparency and aid in resolving the reproducibility crises seen in many fields (Peng, Dominici & Zeger, 2006; Peng, 2011; Parker et al., 2016). Second, digitising programs do not allow the integration of metadata at the time of data extraction, such as experimental group or variable names, and sample sizes. This makes the downstream calculations laborious, as information has to be added later using different software. Third, existing programs do not import sets of images for the user to systematically work through. Instead they require the user to manually import images one by one, and export data into individual files, that need to be imported and edited using different software. Finally, digitising programs typically only provide the user with calibrated x,y coordinates from imported figures, and do not differentiate between common plot types that are used to present data. Consequently a large amount of additional data manipulation is required, that is

different across plots types. For example, data are frequently presented in plots with means and standard errors or confidence intervals (Figure 1A), from which the user wants a mean and SD for each group presented. From x,y coordinates, users must manually discern between mean and error coordinates and assign points to groups. Error then needs to be calculated as the deviation from the mean, and then transformed to SD, according to the type of error presented.

Data extraction from figures is therefore an incredibly time-consuming process as existing software does not provide an optimized, reproducible research pipeline to facilitate data extraction and editing. Here, we present an interactive R package, **metaDigitise** (available at <https://github.com/daniel1noble/metaDigitise>), which is designed for large scale, reproducible data extraction from figures, specifically catering to the the needs of meta-analysts. To this end, we provide tools to extract data from common plot types (mean/error plots, box plots, scatter plots and histograms, see Figure 1). **metaDigitise** operates within the R environment making data extraction, analysis and export more streamlined. The necessary calculations are carried out on calibrated data immediately after extraction so that comparable summary statistics can be obtained quickly. Summary data from multiple figures is returned into a single data frame which can be can easily exported or use in downstream analysis within R. Completed digitisations are automatically saved for each figure, meaning users can redraw their digitisations on figures, make corrections and access calibration and proceeded data. This makes sharing figure digitisation and reproducing the work of others simple and easy, and allows meta-analyses to be updated more efficiently.

2 metaDigitise and Reproducibility

The **metaDigitise** package has one main function, `metaDigitise()`, which interactively takes the user through the process of extracting data from figures. `metaDigitise()`

works on a directory containing images of figures copied from primary literature, in .png, .jpg, .tiff, .pdf format, specified to `metaDigitise()` through the `dir` argument. `metaDigitise()` recognizes all the images in the given directory and automatically imports them one by one, allowing the user to extract the relevant information about a figure as they go. It is worth the user thinking carefully about their directory structure early on in their project. Although different directory structures may be used, we would recommend having all figures for one project in a single directory with an informative and unambiguous naming scheme (e.g. `paper_figure_trait.png`). This expedites digitisation by preventing users from having to constantly change directories and / or open new images.

The data from each completed image is automatically saved as a `metaDigitise` object in a separate .RDS file to a `caldat` directory that is created within the parent directory when first executing `metaDigitise()`. These files enable re-plotting and editing of images at a later point (see below). When run, `metaDigitise()` also identifies the images within a directory that have been previously digitised and only imports new images to process. The data of all images is then automatically integrated into the final output. This means that all figures do not need to be extracted at one time and new figures can be added to the directory as the project develops.

This directory structure allows the complete digitisation process to be reproduced at a later stage, shared with collaborators and presented as supplementary materials for a publication. As long as all the images and the `caldat` directory are still in one directory, `metaDigitise()` will be able to reproduce all figure extractions, regardless of the computer it is run on. For an analysis to be updated, new figures can simply be added to the directory and `metaDigitise()` run to incorporate the new data.

3 Image Processing

Running `metaDigitise()` presents the user with three options; ‘Process new images’, ‘Import existing data’ or ‘Edit existing data’. Selecting ‘Process New Images’ starts the digitisation process on images within the directory that have not previously been digitised; the other functions are discussed below.

For all plot types, `metaDigitise()` requires the user to calibrate the axes in the figure, by clicking on two known points on the axis in question, and entering the value of those points (Figure 1). `metaDigitise()` then calculates the value of any clicked points in terms of the figure axes. This is based on the calibration used in the **digitize** R package (Poisot, 2011). For mean/error and box plots, only the y-axis is calibrated (Figure 1A,B), assuming the x-axis is redundant. For scatter plots and histograms both axes are calibrated (Figure 1C,D).

As figures may have been copied from older, scanned publications, they may not be perfectly orientated. This makes calibration of the points in the figure problematic. `metaDigitise()` allows users to rotate the image (Figure 2A,B). Furthermore, mean/error plots, box plots and histograms, may be presented with horizontal bars. `metaDigitise()` assumes that bars are vertical, but allows the user to flip the image to make the bars are vertical (Figure 2C,D).

metaDigitise recognises four main types of plot; Mean/error plots, box plots, scatter plots and histograms (1). All plot types can be extracted in a single call of `metaDigitise()` and integrated into one output. Alternatively, users can process different plot types separately, using separate directories. All four plot types are extracted slightly differently (outlined below). Upon completing all images, or quitting, either summarised or calibrated data is returned (specified by the user through the `summary` argument). Summarised data consists of a mean, SD and sample size, for each identified group within the plot (should multiple groups exist). In the case of scatter

132 plots, the correlation coefficient between x and y variables within each identified group
133 is also returned. Calibrated data consists of a list with slots for each of the four figure
134 types, containing the calibrated points that the user has clicked. This may be
135 particularly useful in the case of scatter plots.

136 **3.1 Mean/Error and Box Plots**

137 `metaDigitise()` handles mean/error and box plots in a very similar way. For each
138 mean/box, the user enters group names and sample sizes. If the user does not enter a
139 sample size at the time of data extraction (if, for example, the information is not readily
140 available) a SD is not calculated. Sample sizes can, however, be entered at a later time
141 (see below). For mean/error plots, the user clicks on an error bar and the mean. Error
142 bars above or below the mean can be clicked, as sometimes one is clearer than the
143 other. `metaDigitise()` assumes that the error bars are symmetrical. Points are
144 displayed where the user has clicked, with the error in a different colour to the mean
145 (Figure 1A). The user also enters the type of error used in the figure: SD, standard
146 error (SE) or 95% confidence intervals (CI95). For box plots, the user clicks on the
147 maximum, upper quartile, median, lower quartile and minimum. For both plot types,
148 the user can add, edit or remove groups. Three functions, `error_to_sd()`,
149 `rqm_to_mean()` and `rqm_to_sd()`, that convert different error types to SD, box plot data
150 to mean and box plot data SD, respectively, are also available in the package (see
151 supplements for further details of these conversions).

152 **3.2 Scatter plots**

153 Users can extract points from multiple groups from scatter plots. Different groups are
154 plotted in different colours and shapes to enable them to be distinguished, with a legend
155 at the bottom of the figure (Figure 1C). Mean, SD and sample size are calculated from

the clicked points, for each group. Data points may overlap with each other making it impossible to know whether points have been missed. This may result in the sample size of digitised groups conflicting with what is reported in the paper. For example, in Figure 1C only 49 points have been clicked when the sample size is known to be 50. Hence, **metaDigitise** also provides the user with the option to input known sample sizes directly. Nonetheless, it is important to recognise the impact that overlapping points can have on summary statistics, and in particular on sampling variance.

3.3 Histograms

The user clicks on the top corners of each bar, which are drawn in alternating colours (Figure 1D). Bars are numbered to allow the the user to edit them. As with scatter plots, if the sample size from the extracted data does not match a known sample size, the user can enter an alternate sample size. The calculation of mean, SD and sample size from this data is shown in the supplements.

4 Importing and Editing Previously Digitised data

metaDigitise is also able to re-import, edit and re-plot previously digitised figures. When running `metaDigitise()`, the user can choose to ‘Import existing data’, which returns previously digitised data, from single or all figures. Alternately, the `getExtracted()` function returns the data of previous digitisations, but without user interaction, allowing easier integration into larger scripts. ‘Edit existing data’ allows the user to re-plot or edit information or digitisations that have previously been done.

177 4.1 Adding Sample Sizes to Previous Digitisations

178 In many cases sample sizes may not be readily available when digitising figures. This
179 information does not need to be added at the time of digitisation. To expedite finding
180 and adding these sample sizes at a later point, **metaDigitise()** has a specific edit
181 option that allows users to enter previously omitted sample sizes. This first identifies
182 missing sample sizes in the digitised output, re-plots the relevant figures and prompts
183 the user to enter the sample sizes for the relevant groups in the figure, one by one.

184 5 Software Validation

185 In order to evaluate the consistency of digitisation with **metaDigitise** between users,
186 fourteen people digitized identical sets of 14 images created from a simulated dataset
187 (see supplements). We found no evidence for any inter-observer variability in
188 digitisations for the mean (ICC = 0, 95% CI = 0 to 0.029, $p = 1$), SD (ICC = 0, 95%
189 CI = 0 to 0.033, $p = 0.5$) or correlation coefficient (ICC = 0.053, 95% CI = 0 to 0.296,
190 $p = 0.377$). There was little bias between digitised and true values, on average 1.63%
191 (mean = 0.02%, SD = 4.9%, $r = -0.03\%$) and there were small absolute differences
192 between digitised and true values, on average 2.18% (mean = 0.40%, SD = 5.81%, $r =$
193 0.33%) across all three summary statistics. SD estimates from digitisations are clearly
194 most error prone. The mean absolute differences for each plot type clearly show that
195 this effect is driven by extraction from box plots and histograms (% difference; box plot:
196 15.805, histogram: 5.210, mean/error: 1.500, scatter plot: 0.433). SD estimation from
197 box plot summary statistics is known to be more error prone, especially at small sample
198 sizes (Wan et al., 2014).

199 We also used simulated data to test the accuracy of digitisations with respect to known
200 values (see supplements). **metaDigitise** was extremely accurate at matching clicked
201 points to their true values essentially being perfectly correlated with the true simulated

202 data for both the x -variable (Pearson's correlation; $r \hat{=} 0.999$, $t = 2137.4$, $df = 78$,
203 $p < 0.001$) and y -variable ($r \hat{=} 0.999$, $t = 1897.8$, $df = 78$, $p < 0.001$) in
204 scatterplots.

205 6 Limitations

206 Although **metaDigitise** is very flexible and provides functionality not seen in any other
207 package, there are some functions that it does not perform (see Table S1). Notably
208 **metaDigitise** lacks automated point detection. However, from our experience, manual
209 digitising is more reliable and often equally as fast. Given the variation in image
210 quality, calibration for automatic point detection needs to be done for each figure
211 individually. Additionally, auto-detection often misses points which then need to be
212 manually added. Based on tests of **metaDigitise** (see above), figures can be extracted in
213 around 1-2 minutes, including the entry of metadata. As a result, we do not believe
214 that current automated point detection techniques provide substantial benefits in terms
215 of time or accuracy.

216 **metaDigitise** also (currently) lacks the ability to zoom in on figures. Zooming may
217 enable users to gain greater accuracy when clicking on points. However, from our own
218 experience (see results above), with a reasonably sized screen accuracy is already high,
219 and so relatively little gain is to be had from zooming in on points.

220 In contrast to some other packages **metaDigitise** does not extract lines from figures.
221 Line extraction is not particularly useful for most comparative or meta-analytic work,
222 although we recognise that it may be useful in fields other than meta-analysis. Should a
223 user like to extract lines with **metaDigitise**, we would recommend extracting data as a
224 scatter plot, and clicking along the line in question. A model can then be fitted to these
225 points (accessed by choosing to return calibrated rather than summary data) to
226 estimate the parameters needed.

227 7 Conclusions

228 Increasing the reproducibility of figure extraction for meta-analysis and making this
229 laborious process more streamlined, flexible and integrated with existing statistical
230 software will go a long way in facilitating the production of high quality meta-analytic
231 studies that can be updated in the future. We believe that **metaDigitise** will improve
232 this research synthesis pipeline, and will hopefully become an integral package that can
233 be added to the meta-analysts toolkit.

234 Acknowledgments

235 We thank the I-DEEL group and colleagues at UNSW for for testing, providing
236 feedback and digitising including: Rose O’Dea, Fonti Kar, Malgorzata Lagisz, Julia
237 Riley, Diego Barneche, Erin Macartney, Ivan Beltran, Gihan Samarasinghe, Dax Kellie,
238 Jonathan Noble, Yian Noble and Alison Pick. J.L.P. was supported by a Swiss National
239 Science Foundation Early Mobility grant (P2ZHP3_164962), D.W.A.N. was supported
240 by an Australian Research Council Discovery Early Career Research Award
241 (DE150101774) and UNSW Vice Chancellors Fellowship and S.N. an Australian
242 Research Council Future Fellowship (FT130100268).

243 Author Contributions

244 J.L.P. and D.W.A.N. conceived the study and J.L.P., S.N. and D.W.A.N. developed the
245 idea. J.L.P. and D.W.A.N. developed the R-package. J.L.P. and D.W.A.N. wrote the
246 first draft of the paper and J.L.P., S.N. and D.W.A.N. contributed substantially to
247 subsequent revisions of the manuscript and gave final approval for publication.

248 References

- 249 Arizona-Software (2008) *GraphClick Software, Version 3.0*.
- 250 Parker, T.H., Forstmeier, W., Koricheva, J., Fidler, F., Hadfield, J., En Chee, Y., Kelly,
251 C.D., Gurevitch, J. & Nakagawa, S. (2016) Transparency in Ecology and Evolution:
252 Real Problems, Real Solutions. *Trends in Ecology and Evolution*, **31**, 711–719.
- 253 Peng, R.D. (2011) Reproducible research in computational science. *Science*, **334**, 1226.
- 254 Peng, R.D., Dominici, F. & Zeger, S.L. (2006) Reproducible epidemiologic research.
255 *American Journal of Epidemiology*, **163**, 783–789.
- 256 Poisot, T. (2011) The digitize package: extracting numerical data from scatterplots.
257 *The R Journal*, **3**, 25–26.
- 258 Rohatgi, A. (2017) *WebPlotDigitizer Software, Version 4.0*. Austin, Texas, USA.
- 259 Tummers, B. (2006) *DataThief Software, Version 3.0*.
- 260 Wan, X., Wang, W., Liu, J. & Tong, T. (2014) Estimating the sample mean and
261 standard deviation from the sample size, median, range and/or interquartile range.
262 *BMC Medical Research Methodology*, **14**, 135.

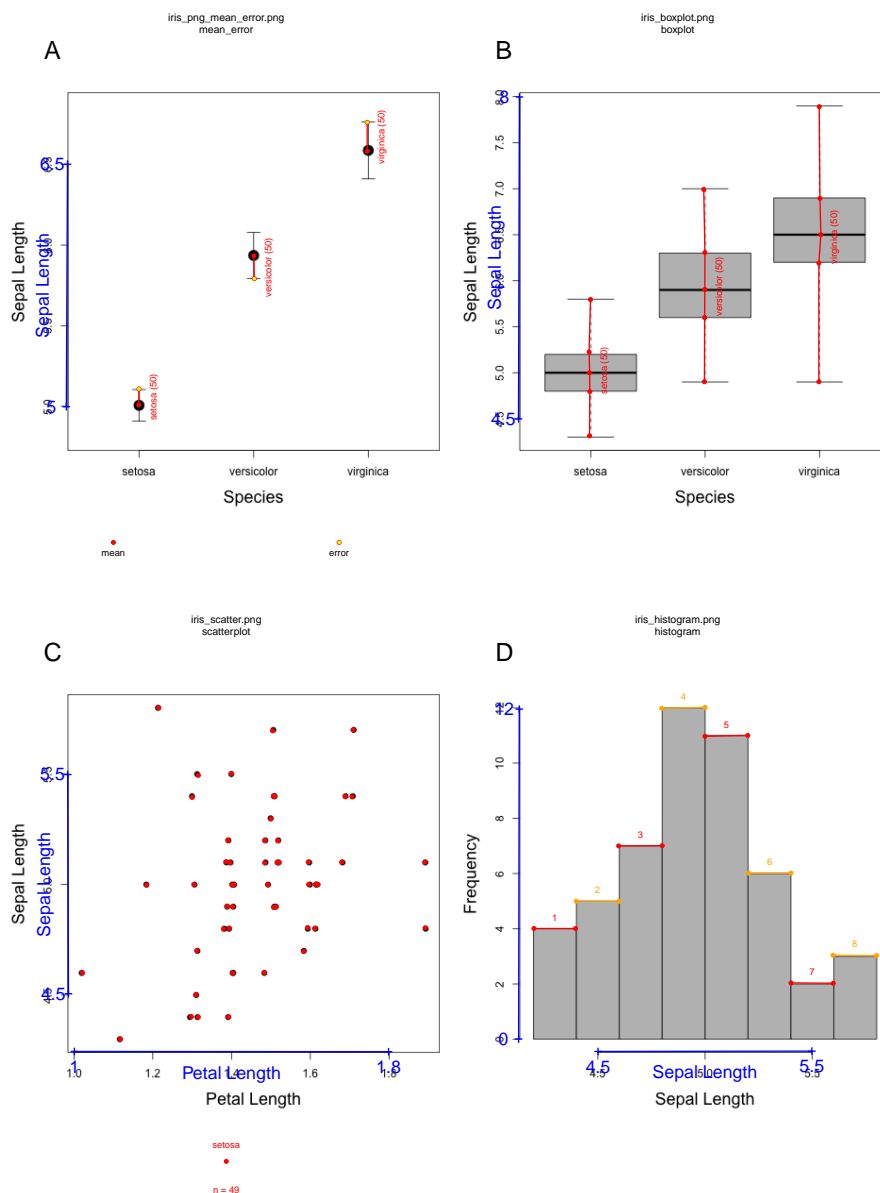


Figure 1: Four plot types that **metaDigitise** is designed to extract data from: A) mean/ error plot, B) box plot, C) scatter plot and D) histogram. Data is taken from the iris dataset in R. A and B are plotted with the whole dataset, C and D are just the data for the species *setosa*. Digitisation of the images is shown. All figures are clearly labelled at the top to remind users of the filename and plot type. This reduces errors throughout the digitisation process. Names of the variables and calibration (in blue) are plotted alongside the digitised points. In A) and B), user entered group names and sample sizes are displayed beside the relevant points. In C) the names and sample sizes for each group are shown below the figure.

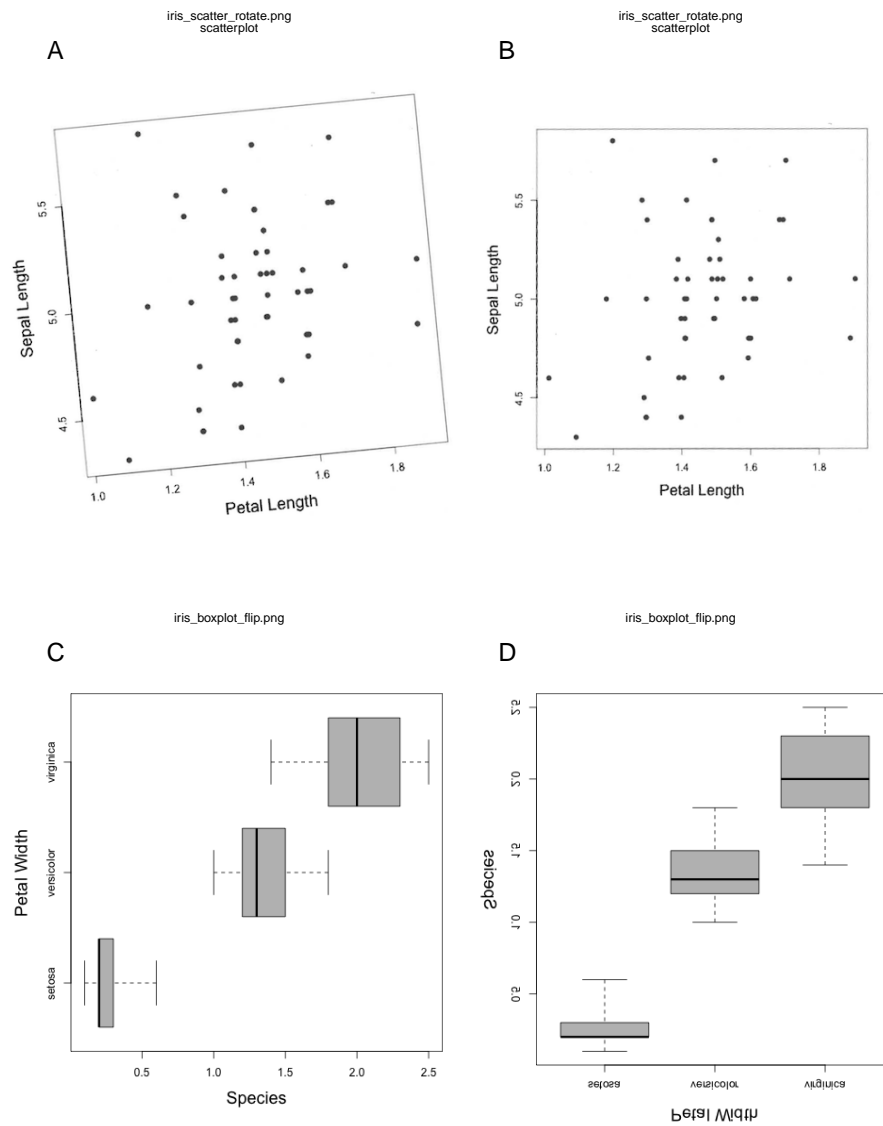


Figure 2: Figure rotation. A) and B) show how non-aligned images can be realigned through user defined rotation. C) and D) show how figures can be re-orientated so as to aid data input.