

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 the fields of ecology and evolution, researchers make use of data presented in primary literature for 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), see Table 1), these tools do not cater to needs of meta-analysts for four main reasons (here we focus on meta-analysis, although many points apply to extraction for 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 and export the resulting digitised data into individual files one-by-one. These data often subsequently 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. The 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 used in downstream analysis within R. Completed digitisations are automatically saved for each figure, meaning users can redraw their digitisations (along with metadata) on figures, make corrections and access calibration and processed (i.e., summarised) data. This makes sharing figure digitisation and reproducing the work of others simple and easy, and allows meta-analyses to be updated more efficiently.

80 2 metaDigitise and Reproducibility

81 The **metaDigitise** package has one main function, `metaDigitise()`, which interactively
82 takes the user through the process of extracting data from figures. `metaDigitise()`
83 works on a directory containing images of figures copied from primary literature, in
84 .png, .jpg, .tiff, .pdf format, specified to `metaDigitise()` through the `dir` argument.
85 `metaDigitise()` recognizes all the images in the given directory and automatically
86 imports them one-by-one, allowing the user to extract the relevant information about a
87 figure as they go. Figures can be organised in different ways for a project, but we would
88 recommend having all figures for one project in a single directory with an informative
89 and unambiguous naming scheme (e.g. `paper_figure_trait.png`). This expedites
90 digitisation by preventing users from having to constantly change directories and / or
91 open new images.

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

100 The complete digitisation process can then be reproduced at a later stage, shared with
101 collaborators and presented as supplementary materials for a publication, regardless of
102 the computer it is run on. For an analysis to be updated, new figures can simply be
103 added to the directory and `metaDigitise()` run to incorporate the new data.

104 3 Image Processing

105 Running `metaDigitise()` presents the user with three options; ‘Process new images’,
106 ‘Import existing data’ or ‘Edit existing data’, which can be used during and after
107 digitisation to execute a range of functions (see Figure 1 – ‘Editing’ and ‘Importing’ are
108 discussed in the next section). Selecting ‘Process New Images’ starts the digitisation
109 process on images within the directory that have not previously been digitised. For all
110 plot types, `metaDigitise()` requires the user to calibrate the axes in the figure, by
111 clicking on two known points on the axis in question, and entering the value of those
112 points (Figure 1). `metaDigitise()` then calculates the value of any clicked points in
113 terms of the figure axes. This is based on the calibration used in the **digitize** R package
114 (Poisot, 2011). For mean/error and box plots, only the y-axis is calibrated (Figure 1),
115 assuming the x-axis is redundant. For scatter plots and histograms both axes are
116 calibrated (Figure 1).

117 As figures may have been copied from older, scanned publications, they may not be
118 perfectly orientated. This makes calibration of the points in the figure problematic.
119 `metaDigitise()` allows users to rotate the image (Figure S2A,B). Furthermore,
120 mean/error plots, box plots and histograms, may be presented with horizontal bars.
121 `metaDigitise()` assumes that bars are vertical, but allows the user to flip the image to
122 make the bars are vertical (Figure S2C,D).

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

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

135 **3.1 Mean/Error and Box Plots**

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

151 **3.2 Histograms**

152 The user clicks on the top corners of each bar, which are drawn in alternating colours
153 (Figure 1C). Bars are numbered to allow the the user to edit them. As with scatter
154 plots, if the sample size from the extracted data does not match a known sample size,

155 the user can enter an alternate sample size. The formulas for calculation of mean, SD
156 and sample size are provided in the supplement.

157 **3.3 Scatter plots**

158 Users can extract points from multiple groups from scatter plots. Different groups are
159 plotted in different colours and shapes to enable them to be distinguished, with a legend
160 at the bottom of the figure (Figure 1D). Mean, SD and sample size are calculated from
161 the clicked points, for each group. Data points may overlap with each other making it
162 impossible to know whether points have been missed. This may result in the sample
163 size of digitised groups conflicting with what is reported in the paper. However, users
164 also have the option to input known sample sizes directly, if required. Nonetheless, it is
165 important to recognise the impact that overlapping points can have on summary
166 statistics, and in particular on sampling variance.

167 **4 Importing and Editing Previously Digitised** 168 **data**

169 **metaDigitise** is also able to re-import, edit and re-plot previously digitised figures.
170 When running **metaDigitise()**, the user can choose to 'Import existing data', which
171 returns previously digitised data, from a single figure or all figures. Alternately, the
172 **getExtracted()** function returns the data of previous digitisations, but without user
173 interaction, allowing easier integration into larger scripts. 'Edit existing data' allows the
174 user to re-plot or edit information for digitisations that have previously be done.

175 4.1 Adding Sample Sizes to Previous Digitisations

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

182 5 Software Validation

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

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

200 data for both the x -variable (Pearson's correlation; $r \geq 0.999$, $t = 2137.4$, $df = 78$,
201 $p < 0.001$) and y -variable ($r \geq 0.999$, $t = 1897.8$, $df = 78$, $p < 0.001$) in
202 scatterplots.

203 6 Limitations

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

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

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

225 7 Conclusions

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

232 Acknowledgments

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

241 Author Contributions

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

246 References

- 247 Arizona-Software (2008) *GraphClick Software, Version 3.0*.
- 248 Bormann, I. (2012) *Digitizelt Software, Version 2.0*. Braunschweig, Germany.
- 249 Lajeunesse, M.J. (2016) Facilitating systematic reviews, data extraction, and
250 meta-analysis with the metagear package for R. *Methods in Ecology and Evolution*, **7**,
251 323–330.
- 252 Parker, T.H., Forstmeier, W., Koricheva, J., Fidler, F., Hadfield, J., En Chee, Y., Kelly,
253 C.D., Gurevitch, J. & Nakagawa, S. (2016) Transparency in Ecology and Evolution:
254 Real Problems, Real Solutions. *Trends in Ecology and Evolution*, **31**, 711–719.
- 255 Peng, R.D. (2011) Reproducible research in computational science. *Science*, **334**, 1226.
- 256 Peng, R.D., Dominici, F. & Zeger, S.L. (2006) Reproducible epidemiologic research.
257 *American Journal of Epidemiology*, **163**, 783–789.
- 258 Poisot, T. (2011) The digitize package: extracting numerical data from scatterplots.
259 *The R Journal*, **3**, 25–26.
- 260 Rohatgi, A. (2017) *WebPlotDigitizer Software, Version 4.0*. Austin, Texas, USA.
- 261 Tummers, B. (2006) *DataThief Software, Version 3.0*.
- 262 Wan, X., Wang, W., Liu, J. & Tong, T. (2014) Estimating the sample mean and
263 standard deviation from the sample size, median, range and/or interquartile range.
264 *BMC Medical Research Methodology*, **14**, 135.

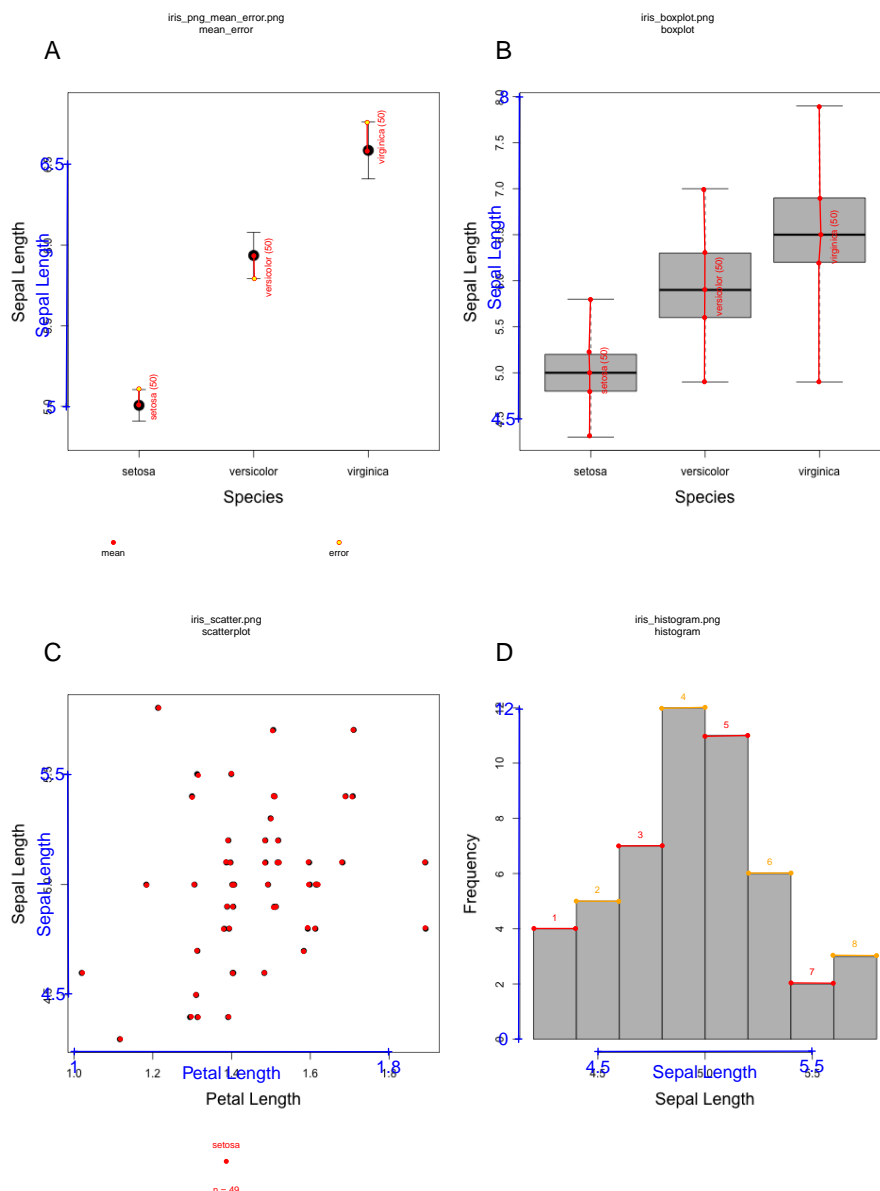


Figure 1: Four plot types that **metaDigitise** is designed to extract data from: A) mean/ error plot, B) box plot, C) scatter plot and D) histogram. Data is taken from the iris dataset in R. A and B are plotted with the whole dataset, C and D are just the data for the species *setosa*. Digitisation 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.

Function	metaDigitise	GraphClick ¹	DataThief ²	DigitizeIt ³	WebPlotDigitizer ⁴	metagear ⁵	digitize ⁶
Scatterplots	✓	✓	✓	✓	✓	✓ ⁷	✓
Mean/error plots	✓	✓	✓	×	×	✓ ⁷	×
Boxplots	✓	×	×	×	×	×	×
Histograms	✓	×	×	×	✓ ⁷	×	×
Graph rotation ⁸	✓	✓	✓	✓	✓	×	×
Grouped Data	✓	✓	×	✓	✓	×	×
Entry of metadata	✓	×	×	×	×	×	×
Summarising data	✓	×	×	×	×	×	×
Multiple image processing	✓	×	×	×	×	×	×
Reproducible ⁹	✓	✓	✓	×	✓	✓	×
Automated point detection	×	✓	×	✓	✓	✓	×
Line extraction	×	✓	✓	✓	✓	×	×
Zoom	×	✓	✓	✓	✓	×	×
Log axis	✓	✓	✓	✓	✓	×	×
Dates	×	×	✓	×	✓	×	×
Asymmetric error bars	×	×	✓	×	×	×	×
Freeware	✓ ¹⁰	✓ ¹¹	✓ ¹¹	×	✓ ¹¹	✓ ¹⁰	✓ ¹⁰

¹ Arizona-Software (2008) ² Tummers (2006) ³ Bornann (2012) ⁴ Rohatgi (2017) ⁵ Lajeunesse (2016) ⁶ Poisot (2011)

⁷ Only automated, no manual extraction.

⁸ Or handles rotated graphs.

⁹ Allows saving, re-plotting and editing of data extraction.

¹⁰ R package.

¹¹ Standalone software.

Table 1: Comparison of functionality between different digitisation softwares.