

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

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

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

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

Abstract

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

27 **1 Introduction**

28 In many different contexts, researchers need to make use of data presented in primary
29 literature. Most notably, this includes meta-analysis, which is becoming increasingly
30 common in many research fields. Meta-analysis uses effect size estimates and their
31 sampling variance, taken from many studies, to understand whether particular effects
32 are common across studies and to explain variation among these effects (Glass, 1976;
33 Borenstein et al., 2009; Koricheva, Gurevitch & Mengersen, 2013; Nakagawa et al.,
34 2017). Meta-analysis relies on descriptive statistics (e.g. means, standard deviations,
35 sample sizes, correlation coefficients) extracted from primary literature that have been
36 reported in the text or tables of research papers. Descriptive statistics are also, however,
37 frequently presented in figures and so need to be manually extracted using digitising
38 programs. While inferential statistics (e.g., t - and F -statistics) are often presented
39 along side descriptive statistics, and can be used to derive effect sizes, descriptive
40 statistics are much more appropriate to use because sources of non-independence in
41 experimental designs can be dealt with more easily (Noble et al., 2017).

42 Although there are several existing tools to perform tasks like this (e.g. **DataThief**
43 (Tummers, 2006), **GraphClick** (Arizona-Software, 2008), **WebPlotDigitizer** (Rohatgi,
44 2017)), these tools are not appropriate for the needs of meta-analysis for three main
45 reasons. First, they typically only provide the user with calibrated x,y coordinates from
46 imported figures, and do not differentiate between common plot types that are used to
47 present data. This means that a large amount of downstream data manipulation is
48 required, that is different across plots types. For example, data are frequently presented
49 in mean/error plots (Figure 1A), for which the user wants a mean and error estimate
50 for each group presented in the figure. With existing programs, from x,y coordinates

51 users must manually discern between mean and error coordinates and assign points to
52 groups. Error then needs to be calculated as the deviation from the mean, and then
53 transformed to a standard deviation, depending on the type of error presented. Second,
54 digitising programs do not easily allow the integration of metadata at the time of data
55 extraction, such as experimental group or variable names, and sample sizes. This makes
56 the downstream calculations more laborious, as the information has to be added later,
57 in most cases using different software. Finally, existing programs do not automatically
58 import new images upon completion to allow the user to systematically work through
59 them. Instead they require the user to manually import images one by one, and export
60 data into individual files, that need to be imported and edited using different
61 software.

62 These present major issues because extracting from figures is an incredibly
63 time-consuming process and existing software does not provide an optimized research
64 pipeline to facilitate data extraction and editing. Furthermore, although meta-analysis
65 is an important tool in consolidating the data from multiple studies, many of the
66 processes involved in data extraction are opaque and difficult to reproduce, making
67 extending studies problematic. Having a tool that facilitates reproducibility in
68 meta-analyses will increase transparency and aid in resolving the reproducibility crises
69 seen in many fields (Peng, Dominici & Zeger, 2006; Peng, 2011; Sandve et al., 2013;
70 Parker et al., 2016; Ihle et al., 2017).

71 Here, we present an interactive R package, **metaDigitise** (available at
72 <https://github.com/daniellnoble/metaDigitise>), which is designed for large scale,
73 reproducible data extraction from figures, specifically catering to the the needs of
74 meta-analysts. To this end, we provide tools specific to data extraction from common
75 plot types (mean/error plots, box plots, scatter plots and histograms, see Figure 1).
76 **metaDigitise** operates within the R environment making data extraction, analysis and
77 export more streamlined. The necessary calculations are carried out on processed data

78 immediately after extraction so that comparable summary statistics can be obtained
79 quickly. Summary data from multiple figures is returned into a single data frame which
80 can be easily exported or use in downstream analysis within R. Completed
81 digitisations are automatically saved for each figure. Users can therefore redraw their
82 digitisations on figures, make corrections and access calibration and proceeded data.
83 This makes sharing figure digitisation and reproducing the work of others simple and
84 easy, and allows meta-analyses to be updated more efficiently.

85 2 Directory Structure and Reproducibility

The **metaDigitise** package was created with the idea that users would have multiple images to extract from and therefore operates in the same way whether the user has one or multiple images. There is 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. The user needs to think carefully about their directory structure early on in their project. Although different directory structures may be used, we would recommend having all files for one project in a single directory with an informative and unambiguous naming scheme for images to help identify the paper and figure that data come from (e.g. *paper_{figure_{trait}}*.png).

86 When **metaDigitise()** is run, it recognizes all the images in a directory and
87 automatically imports them one by one, allowing the user to extract the relevant
88 information about a figure as they go. Having all figures in one directory therefore
89 expedites digitising figures by preventing users from having to constantly change
90 directories and / or open new images. The data from each completed image is
91 automatically saved as a **metaDigitise** object in a separate .RDS file to a **caldat**

92 directory that is created within the parent directory when first executing
93 `metaDigitise()`. These files enable re-plotting and editing of images at a later point
94 (see below).

95 `metaDigitise()` identifies images within a directory that have been previously digitised
96 and only import images that have not been digitised in previous calls of the function.
97 All figures do not, therefore, need to be extracted at one time and new figures can be
98 added to the directory as the project develops. After each image is extracted, the user
99 is asked whether they wish to continue or quit the extraction process. Upon rerunning
100 `metaDigitise()`, previously digitised figures are ignored during processing, but their
101 data is automatically integrated into the final output.

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

108 **3 Image Processing**

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

113 For all plot types, `metaDigitise()` requires the user to calibrate the axes in the figure,
114 by clicking on two known points on the axis in question, and entering the value of those
115 points (Figure 1). `metaDigitise()` then calculates the value of any clicked points in
116 terms of the figure axes. For mean/error and box plots, only the y-axis is calibrated

117 (Figure 1A,B), assuming the x-axis is redundant. For scatter plots and histograms both
118 axes are calibrated (Figure 1C,D).

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

125 **metaDigitise** recognises four main types of plot; Mean/error plots, box plots, scatter
126 plots and histograms, shown in Figure 1. All plot types can be extracted in a single call
127 of `metaDigitise()` and integrated into one output. Alternatively, users can process
128 different plot types separately, using separate directories. All four plot types are
129 extracted slightly differently (outlined below). After completing images, or upon
130 quitting, either summarised or processed data is returned (specified by the user).
131 Summarised data consists of a mean, standard deviation and sample size, for each
132 identified group within the plot (should multiple groups exist). In the case of scatter
133 plots, the correlation coefficient between the points within each identified group is also
134 returned. Alternatively choosing processed data will return a list with slots for each of
135 the four figure types, containing the calibrated points that the user has clicked. This
136 may be particularly useful in the case of scatter plots.

137 3.1 Mean/Error and Box Plots

138 `metaDigitise()` handles mean/error and box plots in a very similar way. For each
139 mean/box, the user is enters group names and sample sizes. If the user does not enter a
140 sample size at the time of data extraction (if, for example, the information is not readily
141 available) a standard deviation (SD) is not calculated. This can, however, be entered at
142 a later time (see below). For mean/error plots, the user clicks on an error bar and the

143 mean. Error bars above or below the mean can be clicked, as sometimes one is clearer
144 than the other. `metaDigitise()` assumes that the error bars are symmetrical. Points
145 are displayed where the user has clicked, with the error in a different colour to the mean
146 (Figure 1A). The user also enters the type of error used in the figure: standard
147 deviation (SD), standard error (SE) or 95% confidence intervals (CI95). For box plots,
148 the user clicks on the maximum, upper quartile, median, lower quartile and minimum.
149 For both plot types, the user can add, edit or remove groups. Three functions,
150 `error_to_sd()`, `rqm_to_mean()` and `rqm_to_sd()`, that convert different error types to
151 SD, box plot data to mean and box plot data SD, respectively, are also available in the
152 package (see SM for further details of these conversions).

153 3.2 Scatter plots

154 Users can extract points from multiple groups from scatter plots. Different groups are
155 plotted in different colours and shapes to enable them to be distinguished, with a legend
156 at the bottom of the figure (Figure 1C). Mean, SD and sample size are calculated from
157 the clicked points, for each group. Data points may overlap with each other making it
158 impossible to know whether points have been missed. However, this will result in the
159 sample size of digitised groups conflicting with what is reported in the paper. For
160 example, in Figure 1C only 49 points have been clicked when the sample size is known
161 to be 50. Hence, **metaDigitise** also provides the user with the option to input known
162 sample sizes directly. Nonetheless, it is important to recognise the impact that
163 overlapping points can have on summary statistics, and in particular on sampling
164 variance.

165 **3.3 Histograms**

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

171 **4 Importing and Editing Previously Digitised** 172 **data**

173 **metaDigitise** is also able to re-import, edit and re-plot previously digitised figures.
174 When running **metaDigitise()**, the user can choose to ‘Import existing data’, which
175 returns previously digitised data. Users can also choose to ‘Edit existing data’ which
176 allows them to re-plot or edit information or digitisations that have previously be done.
177 Points added by mistake can be deleted. The user can add more groups, edit groups
178 (add or remove points) or delete groups and this is automatically re-incorporated in the
179 data.

180 **4.1 Adding Sample Sizes to Previous Digitisations**

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

188 5 Software Validation

189 In order to evaluate the consistency of digitisation using **metaDigitise** between users, we
190 simulated a dataset of two variables with two groups. The same simulated datasets were
191 given to 14 different digitisers to compare the inter-observer variability in digitisations.
192 We also used simulated data to test the accuracy of digitisations with respect to known
193 values (See SM for more details on how simulations were set up).

194 Across the plot types we found no evidence for any inter-observer variability in
195 digitisations for the mean (ICC = 0, 95% CI = 0 to 0.029, $p = 1$), standard deviation
196 (ICC = 0, 95% CI = 0 to 0.033, $p = 0.5$) or correlation coefficient (ICC = 0.053, 95%
197 CI = 0 to 0.296, $p = 0.377$). There were also little bias between digitised and true
198 values, on average 1.63% (mean = 0.02%, SD = 4.9%, $r = -0.03\%$) and overall there
199 were only small absolute differences between digitised and true values, deviating, on
200 average 2.18% (mean = 0.40%, SD = 5.81%, $r = 0.33\%$) across all three summary
201 statistics. SD estimates from digitisations are clearly more prone to error than means or
202 correlation coefficients. If the mean absolute difference is calculated for each plot type,
203 we can see that this effect is driven mainly by extraction from box plots and histograms
204 (% difference; box plot: 15.805, histogram: 5.210, mean/error: 1.500, scatter plot:
205 0.433). SD estimation from box plot summary statistics is known to be more error
206 prone, especially at small sample sizes (Wan et al., 2014).

207 **metaDigitise** was extremely accurate at matching clicked points to their true values
208 essentially being perfectly correlated with the true simulated data for both the
209 x -variable (Pearson's correlation; $r = 0.9999915$, $t = 2137.4$, $df = 78$, $p < 0.001$) and
210 y -variable ($r = 0.9999892$, $t = 1897.8$, $df = 78$, $p < 0.001$) in scatterplots.

211 6 Limitations and Future Extensions

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

222 **metaDigitise** also (currently) lacks the ability to zoom in on plots. Zooming may enable
223 users to gain greater accuracy when clicking on points. However, from our own
224 experience (and indeed from the results above), if you are using a reasonably sized
225 screen then the accuracy is already high, and so there is not much gain to be had from
226 zooming in on points.

227 In contrast to some other packages (Table S1), **metaDigitise** does not extract lines from
228 figures. Line extraction may not particularly useful for most meta-analyses, although we
229 recognise that it may be useful in other fields. Should a user like to extract lines with
230 **metaDigitise**, we would recommend extracting data as a scatter plot, and clicking along
231 the line in question. A model can then be fitted to these points (accessed by choosing to
232 return processed rather than summary data) to estimate the parameters needed.

233 7 Conclusions

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

240 Acknowledgments

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

249 Author Contributions

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

254 References

- 255 Arizona-Software (2008) *GraphClick Software, Version 3.0*.
- 256 Borenstein, M., Hedges, L., Higgins, J. & Rothstein, H. (2009) Introduction to
257 meta-analysis. *John Wiley Sons. Ltd. West Sussex, UK*.
- 258 Glass, G. (1976) Primary, secondary, and meta-analysis research. *Educational*
259 *Researcher*, **5**, 3–8.
- 260 Ihle, M., Winney, I.S., Krystalli, A. & Croucher, M. (2017) Striving for transparent and
261 credible research: practical guidelines for behavioral ecologists. *Behavioral Ecology*,
262 **28**, 348–354.
- 263 Koricheva, J., Gurevitch, J. & Mengersen, K. (2013) Handbook of Meta-Analysis in
264 Ecology and Evolution. *Princeton University Press, Princeton, New Jersey*.
- 265 Nakagawa, S., Noble, D.W., Senior, A.M. & Lagisz, M. (2017) Meta-evaluation of
266 meta-analysis: ten appraisal questions for biologists. *BMC Biology*, **15**, 18; DOI
267 10.1186/s12915-017-0357-7.
- 268 Noble, D.W., Lagisz, M., O’Dea, R.E. & Nakagawa, S. (2017) Nonindependence and
269 sensitivity analyses in ecological and evolutionary meta-analyses. *Molecular Ecology*,
270 **26**, 2410–2425.
- 271 Parker, T.H., Forstmeier, W., Koricheva, J., Fidler, F., Hadfield, J., En Chee, Y., Kelly,
272 C.D., Gurevitch, J. & Nakagawa, S. (2016) Transparency in Ecology and Evolution:
273 Real Problems, Real Solutions. *Trends in Ecology and Evolution*, **31**, 711–719.
- 274 Peng, R.D. (2011) Reproducible research in computational science. *Science*, **334**, 1226.
- 275 Peng, R.D., Dominici, F. & Zeger, S.L. (2006) Reproducible epidemiologic research.
276 *American Journal of Epidemiology*, **163**, 783–789.
- 277 Rohatgi, A. (2017) *WebPlotDigitizer Software, Version 4.0*. Austin, Texas, USA.

278 Sandve, G.K., Nekrutenko, A., Taylor, J. & Hovig, E. (2013) Ten simple rules for
279 reproducible computational research. *PLoS Computational Biology*, **9**, e1003285.

280 Tummers, B. (2006) *DataThief Software, Version 3.0*.

281 Wan, X., Wang, W., Liu, J. & Tong, T. (2014) Estimating the sample mean and
282 standard deviation from the sample size, median, range and/or interquartile range.
283 *BMC Medical Research Methodology*, **14**, 135.

284 **Figures**

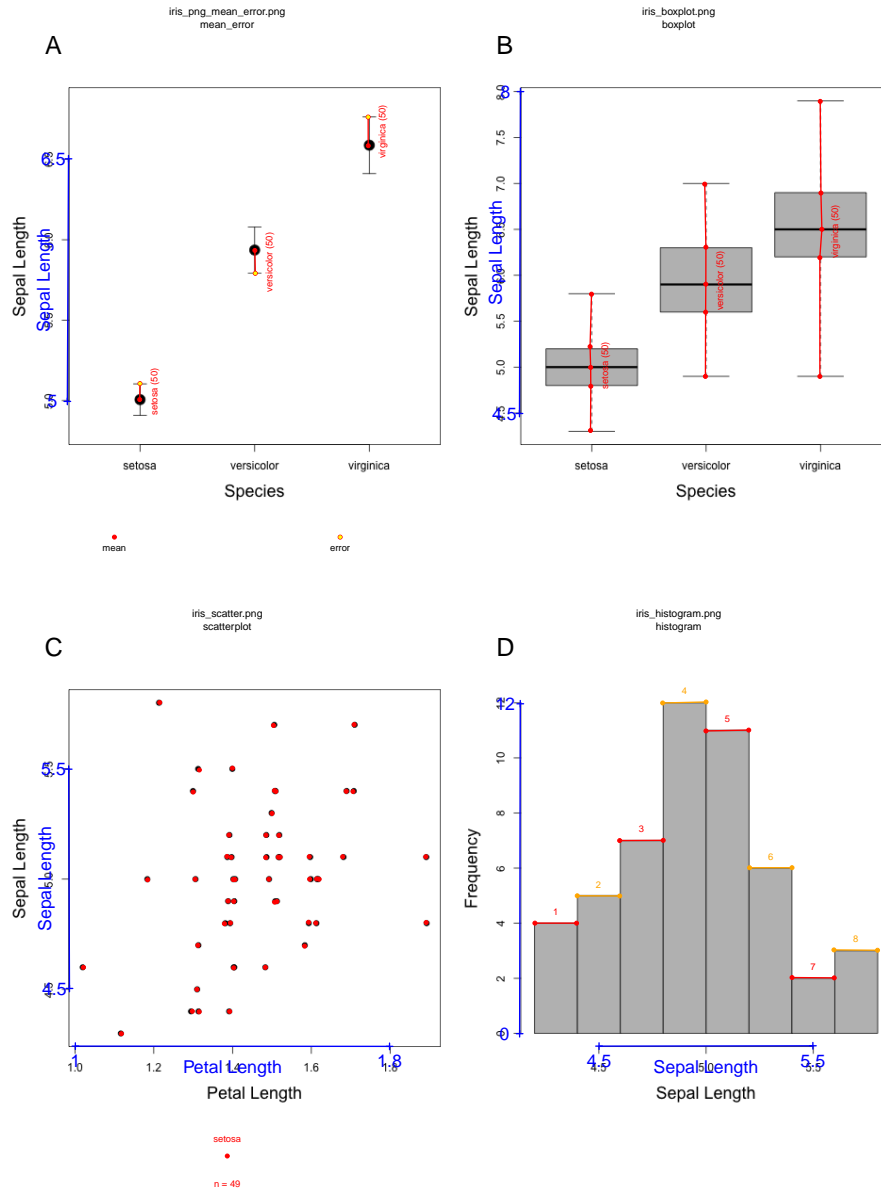


Figure 1: Four plot types that **metaDigitise** is designed to extract data from: A) mean/error plot, B) box plot, C) scatter plot and D) histogram. Data is taken from the iris dataset in R. A and B are plotted with the whole dataset, C and D are just the data for the species *setosa*. Digitisation of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

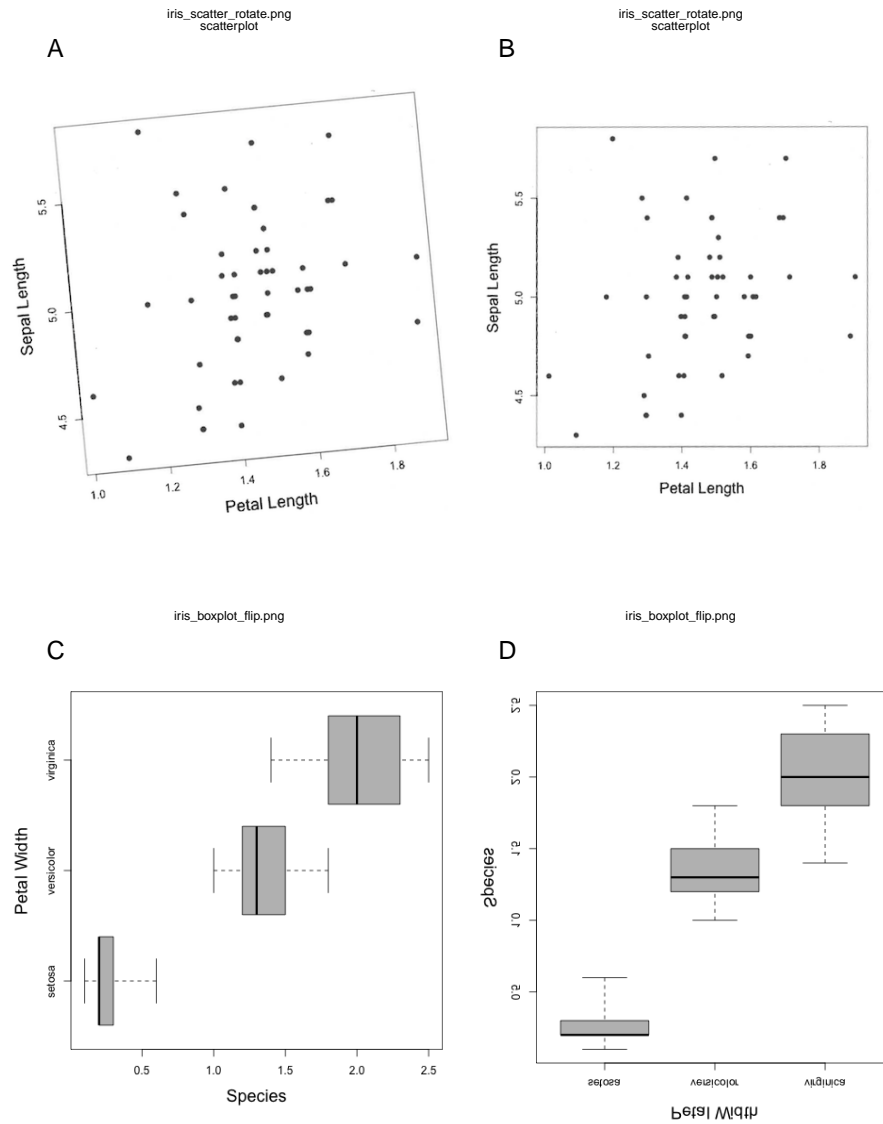


Figure 2: Figure rotation. A) and B) show how non-aligned images can be realigned through user defined rotation. C) and D) show how figures can be re-orientated so as to aid data input.