1    nzffdr: an R package to import, clean and update data from the New Zealand Freshwater

2    Fish Database

3

4    Finnbar Lee[a]* (ORCID: https://orcid.org/0000-0002-9219-1486)

5    Nick Young[b] (ORCID: https://orcid.org/0000-0002-5261-4042)

6

7    [a]School of Environment, University of Auckland, Auckland, New Zealand; [b]Centre for

8    eResearch, University of Auckland, Auckland, New Zealand

9    *corresponding author: finnbar.lee@auckland.ac.nz

10

11

12

13

14

15

16

17

18

19

20

21

22

23    nzffdr: an R package to import, clean and update data from the New Zealand Freshwater

24    Fish Database


25    **Abstract**

26    The New Zealand Freshwater Fish Database (NZFFD) is a repository of more than

27    155,000 records of freshwater fish observations from around New Zealand, maintained

28    by the National Institute of Water and Atmospheric Research (NIWA). Records from

29    the NZFFD can be downloaded using a web interface. The statistical computing

30    language R is now widely used for data wrangling, analysis, and visualisation. Here, we

31    present nzffdr, an open source R software package that: i) allows users to query and

32    download data from the New Zealand Freshwater Fish Database directly in R, ii)

33    provides functions to clean imported data, iii) facilitates the addition of information

34    such as species names and Department of Conservation threat classification status, and

35    iv) a workflow for visualising information from the NZFFD. The nzffdr package aims

36    to standardise, simplify, and speed up a workflow likely already used in an *ad hoc*

37    manner by scientists across New Zealand and abroad.

38

39    Keywords: software; open source; NZFFD, New Zealand, freshwater fish,

40    reproducible workflow


41

42

43

**Introduction**

44

45    The New Zealand Freshwater Fish Database (NZFFD) contains over 155,000

46    observations of freshwater fish (plus freshwater shrimp and kōura) from across New

47    Zealand dating back to 1901 (Crow, 2017). The observations typically include

48    information on sampling location, date, time, fishing method, and the organisation that

49    conducted the survey; less frequently information on the number and size of individuals

50    caught is included. The database is a remarkable asset and is widely used to inform

51    academic and governmental research and decision making (Goodman et al., 2014; Joy

52    & Death, 2004). A limitation of the NZFFD is that it lacks some basic variables that

53    individuals need to add each time they analyse NZFFD data; for example, species'

54    common and scientific names are not included (a 6 letter species code is included), nor

55    is any other taxonomic information (e.g. family), threat classification status, and

56    whether the species is native or introduced. Adding this information each time data are

57    downloaded is not trivial and can be time consuming if many records are downloaded.

58         The statistical computing language R (R Core Team, 2020) has increased in

59    popularity over the last decade, and is now one of the most common programming

60    languages used by ecologists (Lai et al., 2019). R is typically used for data wrangling,

61    analysis, and visualisation and is a popular tool for interrogating NZFFD data (Jellyman

62    & Harding, 2012; Jones & Closs, 2015; Leathwick et al., 2006).

63         Here we present the nzffdr R package. We describe the features of each of the

64    core functions in the nzffdr package and then illustrate how the functions can be used

65    via an analysis of NZFFD data.

**Methodology**

66

67    The nzffdr package has four core functions and four core datasets. The four core

68    functions: i) import NZFFD data from in R. ii) clean up a variety of spelling

69   inconsistencies and add a new variable "form" which describes the sampling habitat e.g.

70   (river, stream, wetland etc.), iii) add missing information such as, family, genus and

71   species names, common names, Department of Conservation threat classification status

72   (Dunn et al., 2017) and whether the species is native or introduced and, iv) import and

73   attaches associated REC data. The four built-in datasets are: i) a subset of 200 rows

74   from the NZFFD that can be accessed without an internet connection and used for

75   exploratory analysis, ii) the different fishing methods included in the NZFFD; it is

76   possible to search the database using these terms so they are provided for reference, iii)

77   scientific and common names of all species included in the NZFFD; the database can be

78   searched by species name (using scientific or common names) so these are provided for

79   reference, iv) a simplified version of the 1:150k NZ map outline available from Land

80   Information New Zealand (https://data.linz.govt.nz/layer/50258-nz-coastlines-topo-

81   150k/)  to facilitate easy mapping of species' distributions.


82   ***Importing data: nzffdr_import()***

83   The *nzffdr_import()* function is used to search the NZFFD and takes input

84   arguments that align with the search options of the NZFFD web user-interface. There

85   are seven search arguments:

86   (1) `catchment`: this refers to the Catchment number, a 6-digit number unique to the

87   reach of interest. Search using the individual reach number (e.g. `catchment =`

88   "702.500"), or for all rivers in a catchment you can use the wildcard search term (e.g.

89   `catchment =` "702%").

90   (2) `river`: search for a river by name; for example, to get all records for the Clutha

91   (`river =` "Clutha").

92     (3) `Location`: search for river by sampling locality for example, to get all records

93     from Awakino (`location` = "Awakino").

94     (4) `fish_method`: search by fishing method used, for example to get all records

95     where fish were caught using a seine net (`fish_method` = "Other net - Seine"). There

96     are currently 59 different possible options for fishing method, a list of all possible

97     fishing method is available via the function *nzffd_method()*.

98     (5) `species`: search for a particular species. There are currently 75 unique species in

99     the NZFFD, a list of all possible species is available via the function

100     `nzffd_species()`. Searches can be made using either common or scientific names

101     and it is possible to search for multiple species at once. e.g. to search for Black mudfish

102     use `species` = "Black mudfish" or `species` = "*Neochanna diversus*" and to search

103     for Black mudfish and Bluegill bully use `species` = c("Black mudfish", "Bluegill

104     bully").

105     (6) `starts`: starting search date.

106     (7) `ends`: ending search date.

107     Not specifying the arguments will return all possible records. The

108     `nzffdr_import()` function requires an internet connection to query NIWA's

109     database.


110     ***Cleaning imported data: nzffd_clean()***

111     While the data imported from NZFFD is generally does not have many errors there are

112     some small inconsistencies (e.g. spelling of river and place names); the

113     `nzffd_clean()` function aims to fix these errors. The first letter of all words in the

114     columns "catchname" and "locality" are capitalised, and any non-alphanumeric

115     characters are removed. Observations in the "time" column are converted to a

116    standardised 24-hour format and nonsensical values (e.g. "0.677") converted to "NA".

117    The organisation column ("org") is converted to all lowercase and has non-

118    alphanumeric characters removed. The NZMS260 map code ("map") is converted to

119    lower case and has any non-three-digit codes converted to "NA". Observations in the

120    catchment name column ("catchname") are standardised, e.g. "Clutha River", "Clutha

121    r" and "Clutha river" all become "Clutha R". Finally, a new variable "form" is added,

122    which defines each observation as one of the following: creek, river, tributary, stream,

123    lake, lagoon, pond, burn, race, dam, estuary, swamp, drain, canal, tarn, wetland,

124    reservoir, brook, spring, gully or NA. The "form" variable is created by matching the

125    above "forms" with the "locality" column; therefore, it reflects the description given by

126    the "locality" variable.


127    *Filling in missing data: nzffd_fill()*

128    Additional useful information can easily be added to the NZFFD dataset. The

129    `nzffd_fill()` function adds columns giving the species' common name

130    ("common_name"), scientific name (genus + species, "sci_name"), "family", "genus",

131    "species", the Department of Conservation threat classification status ("threat_class",

132    [Dunn et al., 2017]) and whether the species is native or introduced ("native").

133    Additionally, if the "map" and "altitude" variables have some "NA" values;

134    `nzffd_fill()` can fill most of these by extracting the relevant data from  The

135    NZMS260 map tiles (https://data.linz.govt.nz/layer/51579-nzms-260-map-sheets) and

136    an 8m digital elevation model (https://data.linz.govt.nz/layer/51768-nz-8m-digital-

137    elevation-model-2012) raster, respectively. This function requires an internet

138    connection to query the 8m DEM.

139     *Adding River Environment Classification data: nzffd_add()*

140     Finally, network topology and environmental information from the River Environment

141     Classification (REC) database (Snelder & Biggs, 2004) can be added to the NZFFD

142     data using `nzffd_add()`. This function takes the NZFFD "nzreach" variable and

143     matches it against the corresponding "NZREACH" variable in the REC database, and

144     imports all the associated REC data, adding 24 new columns to the NZFFD dataset.

145     This function requires an internet connection to query the REC database.


146     **Illustration of nzffdr functionality**

147     To demonstrate the utility of the nzffdr package we imported the entire NZFFD into R,

148     cleaned up the imported data, filled in missing data, and added the REC database. We

149     then highlight the usefulness of some of the new variables that nzffdr has added to the

150     NZFFD dataset. Specifically, we map the distribution of native and introduced species,

151     plot the relative proportion of records across habitat forms for each of the *Galaxias*

152     species, highlighting their respective conservation status, and finally use the REC data

153     to show the distance inland that each of the *Galaxias* species has been found.

154        All analysis was carried out using R v 4.1.0 (R Core Team, 2020). The package

155     dplyr v 1.0.6 (Wickham et al., 2021) was used for data wrangling, ggplot2 v 3.3.3

156     (Wickham, 2016) for visualisation, and nzffdr v 1.0.0 (Lee & Young, 2021) used to

157     access and tidy the NZFFD data. The code used to generate the results presented here is

158     available via Figshare (https://doi.org/10.17608/k6.auckland.14776770.v1).

159


160     **Results and discussion**

161     We plotted the distribution of introduced and native species records from the NZFFD

162     (Fig. 1), where the introduced/native variable and the map of New Zealand are provided

163    by the nzffdr R package. We then graphed the relative number of records occurring

164    across 10 habitat forms for each of the *Galaxias* species, including information about

165    each species' threat classification status (Fig. 2). Habitat form, threat classification, and

166    species common names have all been added to NZFFD data via the nzffdr package.

167    Finally, distance to sea (km) at each of the locations of *Galaxias* species in the NZFFD

168    have been observed at was plotted (Fig. 3). The distance to sea variable is added to the

169    NZFFD data from the River Environment Classification database via the nzffdr

170    package. This analysis illustrates some of the functionality offered by the nzffdr

171    package.

172        Here we have presented an overview of the nzffdr open source software

173    package, which streamlines the importing, tidying, and adding of other important

174    variables to the New Zealand Freshwater Fish database in R. This workflow is likely

175    already being undertaken by researchers across New Zealand and overseas in an *ad hoc*

176    manner. The nzffdr package speeds up this process and contributes to a reproducible

177    workflow.

178    **Acknowledgements**

184 **Data availability statement**

185 The release version of the nzffdr software package described here is archived on the

186 Comprehensive R Archive Network (https://cran.r-project.org) and the latest

187 development version can be installed from https://github.com/flee598/nzffdr. The code

188 used to produce the tables and figure in this manuscript is available via Figshare:

189 https://doi.org/10.17608/k6.auckland.14776770.v1.


190 **References**

191 Crow, S. (2017). New Zealand Freshwater Fish Database. Version 1.2. The National

192       Institute of Water and Atmospheric Research (NIWA). Occurrence Dataset

193       https://doi.org/10.15468/ms5iqu.

194 Dunn, N. R., Allibone, R. M., Closs, G., Crow, S., David, B. O., Goodman, J., Griffiths,

195       M. H., Jack, D., Ling, N., & Waters, J. M. (2017). Conservation status of New

196       Zealand freshwater fishes, 2017. Department of Conservation.

197 Goodman, JM., Dunn, N., Ravenscroft, P., Allibone, R., Boubee, J., David, B.,

198       Griffiths, M., Ling, N., Hitchmough, R., & Rolfe, J. (2014). Conservation status

199       of New Zealand freshwater fish, 2013.

200 Jellyman, P. G., & Harding, J. S. (2012). The role of dams in altering freshwater fish

201       communities in New Zealand. New Zealand Journal of Marine and Freshwater

202       Research, 46(4), 475–489.

203 Jones, P. E., & Closs, G. P. (2015). Life history influences the vulnerability of N ew Z

204       ealand galaxiids to invasive salmonids. Freshwater Biology, 60(10), 2127–2141.

205 Joy, M. K., & Death, R. G. (2004). Application of the index of biotic integrity

206       methodology to New Zealand freshwater fish communities. Environmental

207       Management, 34(3), 415–428.

208    Lai, J., Lortie, C. J., Muenchen, R. A., Yang, J., & Ma, K. (2019). Evaluating the

209         popularity of R in ecology. Ecosphere, 10(1), e02567.

210    Leathwick, J. R., Elith, J., & Hastie, T. (2006). Comparative performance of generalized

211         additive models and multivariate adaptive regression splines for statistical

212         modelling of species distributions. Ecological Modelling, 199(2), 188–196.

213    Lee, F., & Young, N. (2021). nzffdr: Import, Clean and Update Data from the New

214         Zealand Freshwater Fish Database. https://CRAN.R-project.org/package=nzffdr

215    R Core Team. (2020). R: A Language and Environment for Statistical Computing. R

216         Foundation for Statistical Computing. https://www.R-project.org/

217    Snelder, T., & Biggs, B. (2004). New Zealand River Environment Classification User

218         Guide. Ministry for the Environment, Auckland.

219    Wickham, H. (2016). ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag.

220         https://ggplot2.tidyverse.org

221    Wickham, H., François, R., Henry, L., & Müller, K. (2021). dplyr: A Grammar of Data

222         Manipulation. https://CRAN.R-project.org/package=dplyr

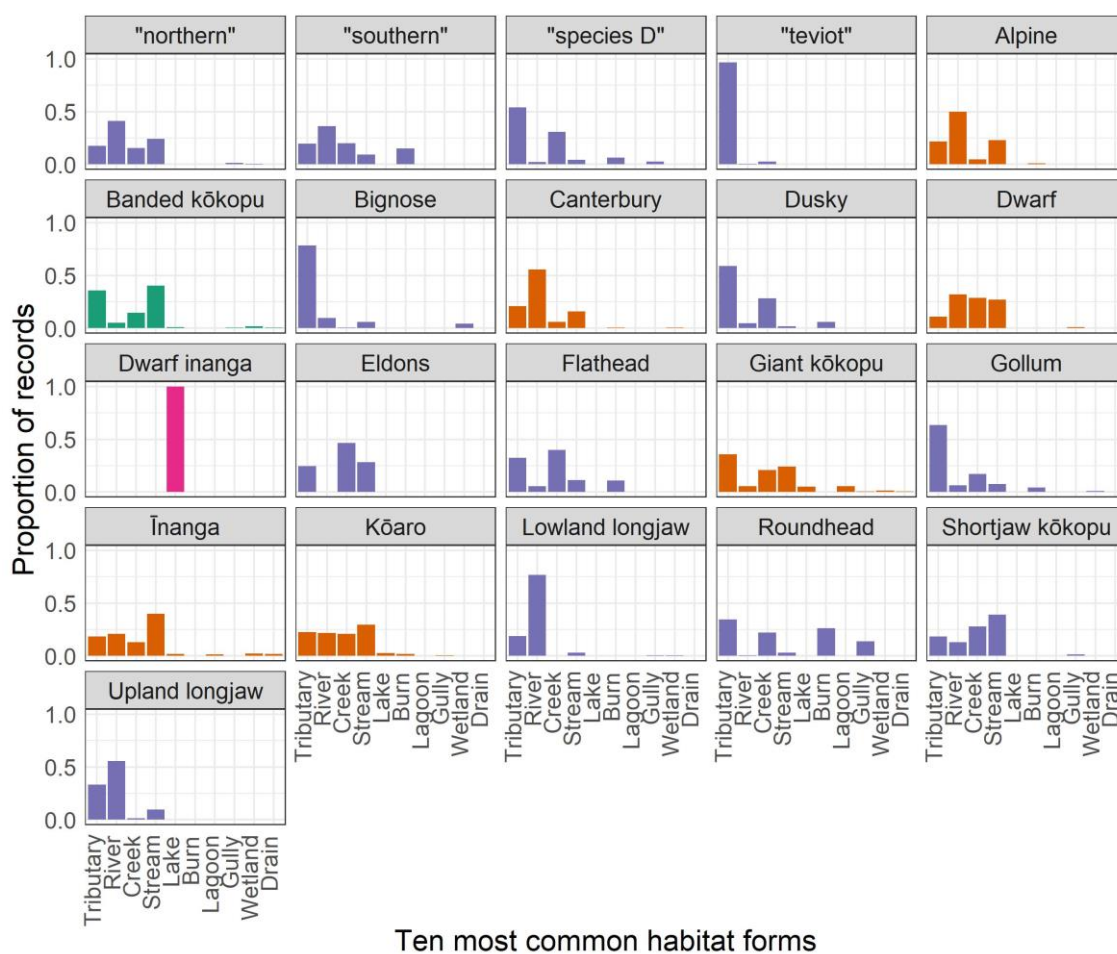223

224

225

226

227

228

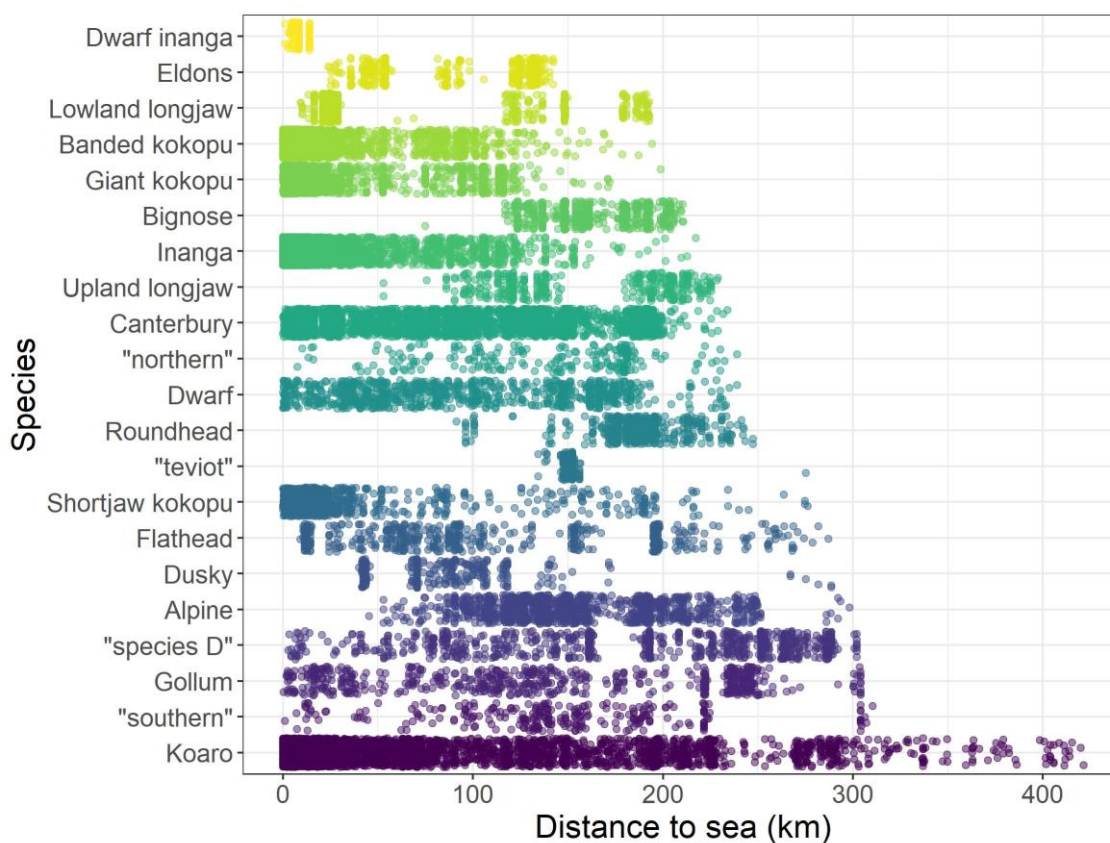229

230

231

232

233

234

235 **Figures**



236   Figure 1. Distribution of introduced and native species records from the NZFFD, where

237   the introduced/native variable and the map of New Zealand are provided by the nzffdr R

238   package.

Figure 2. The relative number of records occurring across 10 habitat forms for each of the *Galaxias* species, the total number of observations for each species is given in parentheses. The habitat form and threat classification variables have been added to NZFFD data via the nzffdr R package.

244

Figure 3. Distance to sea (km) that each of the *Galaxias* species in the NZFFD have

been observed. The distance to sea variable is added to the NZFFD data from the River

Environment Classification database via the nzffdr R package.