

1 **Automated residual plot assessment with the R package `autovi`**
2 **and the Shiny app `autovi.web`**

3 Weihao Li^{1,2}, Dianne Cook¹, Emi Tanaka², Susan VanderPlas³ and Klaus
4 Ackermann¹

5 *Monash University, The Australian National University and University of
6 Nebraska*

Summary

7 Visual assessment of residual plots is a common approach for diagnosing linear
models, but it relies on manual evaluation, which does not scale well and can
lead to inconsistent decisions across analysts. The lineup protocol, which embeds
the observed plot among null plots, can reduce subjectivity but requires even
more human effort. In today's data-driven world, such tasks are well-suited for
automation. We present a new R package that uses a computer vision model to
automate the evaluation of residual plots. An accompanying Shiny app is provided
for ease of use. Given a sample of residuals, the model predicts a visual signal
strength (VSS) and offers supporting information to help analysts assess model fit.

8 *Key words:* initial data analysis; statistical graphics; data visualization; visual inference;
computer vision; machine learning; hypothesis testing; regression analysis;
model diagnostics

9 1. Introduction

10 Regression analysis is a widely used statistical modeling technique for data in many
11 fields. There is a vast array of software for conducting regression modeling and
12 generating diagnostics. The package `lmltest` ([Zeileis & Hothorn 2002](#)) provides a suite
13 of conventional tests, while the `stats` package ([R Core Team 2022](#)) offers standard
14 diagnostic plots such as residuals vs. fitted values, quantile-quantile (Q-Q) plots, and
15 residuals vs. leverage plots. Additional packages like `jtools` ([Long 2022](#)), `olsrr` ([Hebbali
2024](#)), `rockchalk` ([Johnson 2022](#)), and `ggResidpanel` ([Goode & Rey 2019](#)) deliver similar

¹ Department of Econometrics and Business Statistics, Monash University, Wellington Road, VIC 3800, Australia

² Biological Data Science Institute, The Australian National University, 46 Sullivan's Creek Road, ACT 2600, Australia

³ Department of Statistics, University of Nebraska, Hardin Hall, 3310 Holdrege St Suite 340, Lincoln, NE 68583, United States

Email: patrick.li@anu.edu.au

graphical diagnostics, often with enhanced aesthetics or interactive features. These tools collectively produce the core diagnostic plots outlined in the classical text by Cook & Weisberg (1982). The `ecostats` package (Warton 2023) extends these diagnostics by incorporating simulation envelopes into residual plots. Meanwhile, DHARMA (Hartig 2022) compares empirical quantiles (0.25, 0.5, and 0.75) of scaled residuals to their theoretical counterparts, with a strong focus on identifying model violations such as heteroscedasticity, misspecified functional forms, and issues specific to generalized linear and mixed-effect models, like over/under-dispersion. It also provides conventional test annotations to reduce the risk of misinterpretation.

However, relying solely on subjective assessments of these plots can lead to issues such as over-interpreting random patterns as model violations. Li et al. (2024a) demonstrated that visual inference methods, particularly those using the lineup protocol (Buja et al. 2009), offer more practical and reliable assessments of residual patterns than conventional tests, as they are less sensitive to minor departures. Packages such as `nullabor` (Wickham et al. 2020), `HLMdiag` (Loy & Hofmann 2014), and `regressinator` (Reinhart 2024) support this approach by enabling users to compare observed residual plots with plots generated under null hypothesis, thereby helping to quantify the significance of any detected patterns.

As noted in Li et al. (2024b), the lineup protocol has significant limitations in large-scale applications due to its reliance on human labor. To overcome this constraint, a computer vision model was developed alongside a corresponding statistical testing procedure to automate the assessment of residual plots. The model takes as input a residual plot and a set of auxiliary variables (such as the number of observations) and outputs a predicted visual signal strength (VSS). This VSS estimates the degree of deviation between the residual distribution of the fitted model and the reference distribution expected under correct model specification.

To make the statistical testing procedure and trained computer vision model widely accessible, we developed the R package `autovi` along with a companion web interface, `autovi.web`, which allows users to automatically assess their residual plots using the trained computer vision model.

The remainder of this paper is structured as follows: Section 2 introduces the definition and computation of visual signal strength. Section 3 expands on the computation of the null and bootstrapped residuals. Section 4 provides a detailed documentation of the `autovi` package, including its usage and infrastructure. Section 5 focuses on the

51 `autovi.web` interface, describing its design and usage, along with illustrative examples.
 52 Finally, Section 6 presents the main conclusions of this work.

53 **2. Definition and computation of visual signal strength**

54 To train a computer vision model, a measure of the visible pattern in a plot is needed.
 55 We call this the **visual signal strength** (VSS), which measures how prominently a
 56 specific set of visual patterns appears in an image. This can be computed for a training
 57 set of data, and plots, where the generating distributions are specified.
 58 In the context of classical normal linear regression model diagnostics, VSS describes
 59 the clarity of visual patterns on a diagnostic plot that may indicate model violations.
 60 Violations can be categorized as weak, moderate, or strong, but here we treat it
 61 as a continuous positive real variable. Importantly, its interpretation depends on
 62 how it is linked to a function of the data or the underlying data generating process.
 63 Consequently, the calculation of VSS can vary across different model classes or within
 64 the same model, depending on the generating function.
 65 VSS estimates the distance between the residual distribution of a fitted classical normal
 66 linear regression model and a reference distribution (see [Li et al. 2024b](#), for details).
 67 The distance measure is based on the Kullback-Leibler (KL) divergence:

$$D = \log(1 + D_{KL}),$$

68 where D_{KL} is given by:

$$D_{KL} = \int_{\mathbb{R}^n} \log \frac{p(\hat{\mathbf{e}})}{q(\hat{\mathbf{e}})} p(\hat{\mathbf{e}}) d\hat{\mathbf{e}}, \quad (1)$$

69 here, $\hat{\mathbf{e}}$ denotes the residual vector from the regression model, and $p(\cdot)$ and $q(\cdot)$ are
 70 the probability density functions of the reference residual distribution P and the true
 71 residual distribution Q, respectively.

72 This distance measure depends on knowledge of the true residual distribution, which
 73 is unknown in practice. To compute D_{KL} for the training samples, Equation (1)
 74 takes different forms depending on the specific model violations. For instance, where
 75 necessary higher-order predictors, \mathbf{Z} , and their corresponding parameter, β_Z , are
 76 omitted from the fitted linear model, the distance measure can be expanded as follows:

$$D_{KL} = \frac{1}{2} (\boldsymbol{\mu}_z^\top (\text{diag}(\mathbf{R}\sigma^2))^{-1} \boldsymbol{\mu}_z),$$

77 where $\boldsymbol{\mu}_z = \mathbf{R}\mathbf{Z}\boldsymbol{\beta}_z$, $\mathbf{R} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top$ and \mathbf{X} is the design matrix of the
78 regression model.

79 The computer vision model approximates this mapping from a set of residuals to
80 its corresponding distance measure. It is trained on a large number of synthetic
81 regression models, each designed to simulate specific violations of classical linear
82 regression assumptions. These models incorporate non-linearity through Hermite
83 polynomial transformations of predictors, heteroskedasticity by making the error
84 variance a predictor-dependent function, and non-normality by drawing residuals from
85 distributions such as discrete, uniform, and lognormal. Both simple and multiple linear
86 regression structures are used, with controlled parameters to generate diverse and
87 complex residual patterns. Since the data-generating process is known, the distance
88 measure D can be explicitly calculated, enabling supervised training. The computer
89 vision model takes a residual plot as input and outputs the corresponding distance
90 measure, learning to quantify model violations directly from visual patterns. Additional
91 details are provided in [Li et al. \(2024b\)](#).

92 3. Definition and simulation of null and bootstrapped residuals

93 In the subsequent sections, we will frequently refer to null residuals and bootstrapped
94 residuals, so it is helpful to first define and explain how they are generated.

95 **Null residuals** are used to generate null plots within the lineup protocol framework,
96 serving as the foundation for the statistical testing in our automated residual plot
97 assessment. Specifically, they represent residuals generated under the null hypothesis
98 that the model is correctly specified. A common method for simulating null residuals
99 in linear regression involves sampling from a normal distribution with mean zero and
100 variance equal to the estimated variance of the error term. These simulated residuals
101 and their corresponding plots depict what one would expect from a correctly specified
102 model. If the true residual plot exhibits noticeable deviations from these null plots, it
103 may suggest model misspecification.

104 Our computer vision model is trained to assign lower VSS to null plots and higher VSS
105 to plots that display distinct patterns. Accordingly, statistical testing is performed
106 by computing the proportion of null plots whose VSS equals or exceeds that of the

107 observed residual plot. This proportion serves as a p -value for a one-sided hypothesis
108 test.

109 **Bootstrapped residuals** are obtained by refitting the model on bootstrap samples,
110 which are generated by sampling individual observations with replacement from the
111 original dataset. The residual plots obtained from these refitted models are evaluated
112 using the same computer vision model. The predicted VSS from the bootstrapped
113 plots provide an empirical estimate of the variation in the VSS of the observed
114 residual plot. By examining the proportion of bootstrapped plots that also exhibit
115 significant violations, we can assess whether the original conclusion is robust to
116 sampling variability.

117 4. R package: autovi

118 The main purpose of **autovi** is to provide rejection decisions and p -values for testing
119 the null hypothesis (H_0) that the regression model is correctly specified. The package
120 provides automated interpretation of residual plots using computer vision. The name
121 **autovi** stands for **automated visual inference**. This functionality can be accessed
122 through the R package **autovi**, or through a web interface, **autovi.web**, which enables
123 users to perform analyses without installing R, Python, or their associated dependencies
124 locally.

125 4.1. Motivation

126 Figure 1 shows three sets of plots of residuals against fitted values. The simulated
127 example in (a) might be interpreted as a heteroscedastic pattern, however the
128 automated reading would predict this to have a visual signal strength (VSS) of
129 1.53, with a corresponding p -value of 0.25. This means it would be interpreted as
130 a good residual plot, that there is nothing in the data to indicate a violation of
131 model assumptions. Skewness in the predictor variables is generating the apparent
132 heteroscedasticity, where the smaller variance in residuals at larger fitted values is
133 due to smaller sample size only. The Breusch-Pagan test (Breusch & Pagan 1979) for
134 heteroscedasticity would also not reject this as good residual plot.

135 The data in (b) is generated by fitting a linear model predicting `mpg` based on `hp`
136 using the `datasets::mtcars`. It is a small data set, and there is a hint of nonlinear
137 structure not captured by the model. The automated plot reading would predict a
138 VSS of 3.57, which has a p -value less than 0.05. That is, the nonlinear structure is
139 most likely real, and indicates a problem with the model. The conventional test, a

140 Ramsey Regression Equation Specification Error Test (RESET) (Ramsey 1969) would
 141 also strongly detect the nonlinearity.

142 The third example is generated using the `surreal` package (Balamuta 2024), where
 143 structure residuals are embedded in the data. In this case, a quote inspired by Tukey,
 144 “visual summaries focus on unexpected values”, is used to define the residual structure.
 145 The automated plot reading predicts the VSS to be 5.87, with a p -value less than
 146 0.05. Visually, the structure is strikingly clear, but a RESET test for nonlinear
 147 structure would not report a problem. (It would be detected by a Breusch-Pagan for
 148 heteroscedasticity and also Shapiro-Wilk test (Shapiro & Wilk 1965) for non-normality.)

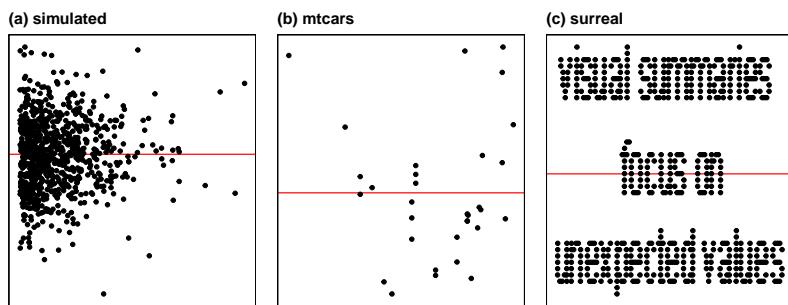


Figure 1. Reading residual plots can be a difficult task, particularly for students new to statistical modeling. The `autovi` package makes it easier. Here are three examples of residual plots, which may appear to have structure. According to `autovi`, the visual signal strengths (VSS) of these three examples are approximately (a) 1.53, (b) 3.57, (c) 5.87, resulting in (b), (c) being significant violations of good residuals, but (a) is consistent with a good residual plot.

149 4.2. Implementation

150 The `autovi` package is built on the `bandicoot` object-oriented programming (OOP)
 151 system (Li 2024), marking a departure from R’s traditional S3 generic system. This
 152 OOP architecture enhances flexibility and modularity, allowing users to redefine key
 153 functions through overriding.

154 The `autovi` infrastructure effectively integrates multiple programming languages and
 155 libraries into a comprehensive analytical tool. It relies on five core libraries from
 156 Python and R, each playing a critical role in the analysis pipeline. In Python, `pillow`
 157 (Clark et al. 2015) handles image processing tasks such as reading and resizing PNG
 158 files of residual plots, then converting them into input tensors for further analysis.
 159 TensorFlow (Abadi et al. 2016), a key component of modern machine learning, is used
 160 to predict the VSS of these plots using a pre-trained convolutional neural network.

161 In the R environment, `autovi` utilizes several libraries. `ggplot2` (Wickham 2016)
 162 generates the initial residual plots, saved as PNG files for visual input. `cassowaryr`
 163 (Mason et al. 2022) computes scagnostics (scatter plot diagnostics), providing numerical
 164 features that capture statistical properties of the plots. These scagnostics complement
 165 the visual analysis by offering quantitative metrics as secondary input to the
 166 computer vision model. `reticulate` (Ushey, Allaire & Tang 2024) enables seamless
 167 communication between R and Python.

168 4.3. Installation

169 The `autovi` package is available on CRAN. It is actively developed and maintained,
 170 with the latest updates accessible on GitHub. This paper uses `autovi` version 0.4.2.
 171 The package includes internal functions to check the current Python environment used
 172 by the `reticulate` package. If the necessary Python packages are not installed in the
 173 Python interpreter, an error will be raised. If you want to select a specific Python
 174 environment, you can do so by calling the `reticulate::use_python()` function before
 175 using the `autovi` package.

176 We recommend using the Shiny app `autovi.web` if users encounter installation
 177 problems.

178 4.4. Usage

179 4.4.1. Numerical summary

180 Three steps are needed to get an automated assessment of a set of residuals and fitted
 181 values:

- 182 1. Load the `autovi` package using the `library()` function.
- 183 2. Create a checker object with a linear regression model.
- 184 3. Call the `check()` method of the checker, which, by default, predicts the VSS for
 185 the true residual plot, 100 null plots, and 100 bootstrapped plots. The method
 186 stores the predictions internally and prints a concise results report.

187 The code to do this is:

```
library(autovi)
checker <- residual_checker(lm(dist ~ speed, data = cars))
checker$check()
```

188 It produces the following summary:

```

189 <AUTO_VI object>
190   Status:
191     - Fitted model: lm
192     - Keras model: (None, 32, 32, 3) + (None, 5) -> (None, 1)
193       - Output node index: 1
194     - Result:
195       - Observed visual signal strength: 3.16 (p-value = 0.0396)
196       - Null visual signal strength: [100 draws]
197         - Mean: 1.274
198         - Quantiles:
199
200           25%    50%    75%    80%    90%    95%    99%
201           0.802  1.111  1.575  1.666  1.919  2.657  3.348
202
203       - Bootstrapped visual signal strength: [100 draws]
204         - Mean: 2.795 (p-value = 0.05941)
205         - Quantiles:
206
207           25%    50%    75%    80%    90%    95%    99%
208           2.455  2.941  3.177  3.300  3.474  3.537  3.668
209
210       - Likelihood ratio: 0.7333 (boot) / 0.06284 (null) = 11.67

```

211 The summary includes observed VSS of the true residual plot and associated p -value
212 of the automated visual test. The p -value is the proportion of null plots (out of the
213 total 100) that have VSS greater than or equal to that of the true residual plot. The
214 report also provides sample quantiles of VSS for null samples and bootstrapped data
215 plots, providing more information about the sampling variability and a likelihood of
216 model violations. The likelihood is computed from the proportion of values greater
217 than the observed VSS in both the bootstrapped data values and the simulated null
218 values.

219 **4.4.2. Visual summary**

220 Users can visually inspect the original residual plot alongside a sample null plot using
221 `plot_pair()` or a lineup of null plot `plot_lineup()`. This visual comparison can
222 clarify why H_0 is either rejected or not, and help identify potential remedies.

```
checker$plot_pair()
```

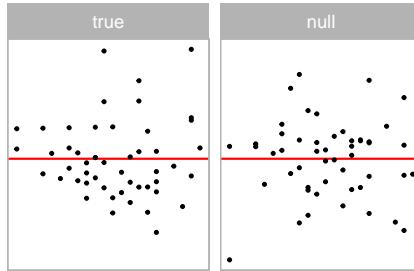


Figure 2. True plot alongside one null plot, for quick comparison.

223 The `plot_pair()` method (Figure 2) displays the true residual plot on the left and a
 224 single null plot on the right. If a full lineup was shown, the true residual plot would
 225 be embedded in a page of null plots. Users should look for any distinct visual patterns
 226 in the true residual plot that are absent in the null plot. Running these functions
 227 multiple times can help any visual suspicions, as each execution generates new random
 228 null plots for comparison.

229 The package offers a straightforward visualization of the assessment result through
 230 the `summary_plot()` function.

```
checker$summary_plot()
```

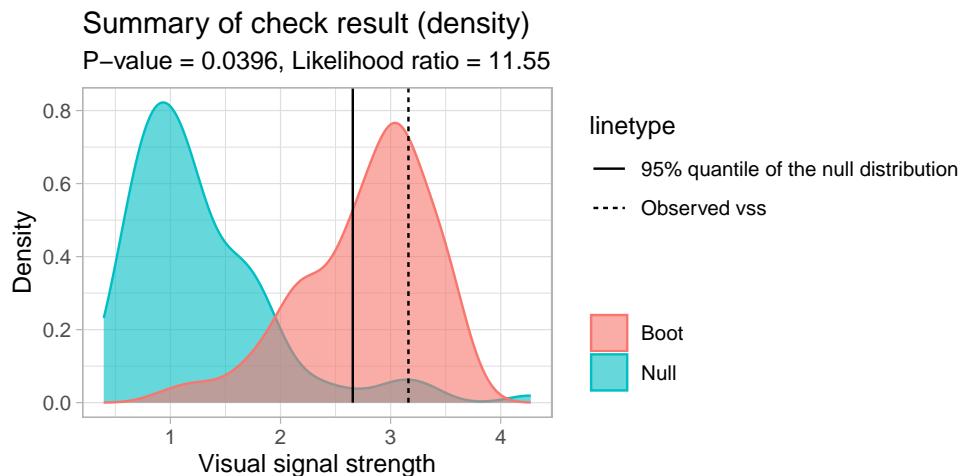


Figure 3. Summary plot comparing the densities of VSS for bootstrapped residual samples (red) relative to VSS for null plots (blue).

231 In the result, shown in Figure 3, the blue area represents the density of VSS for null
 232 residual plots, while the red area shows the density for bootstrapped residual plots.
 233 The dashed line indicates the VSS of the true residual plot, and the solid line marks
 234 the critical value at a 95% significance level. The p -value and the likelihood ratio are
 235 displayed in the subtitle. The likelihood ratio represents the ratio of the likelihood
 236 of observing the VSS of the true residual plot from the bootstrapped distribution
 237 compared to the null distribution.

238 Interpreting the plot involves several key aspects. If the dashed line falls to the right of
 239 the solid line, it suggests rejecting the null hypothesis. The degree of overlap between
 240 the red and blue areas indicates similarity between the true residual plot and null
 241 plots; greater overlap suggests more similarity. Lastly, the portion of the red area to
 242 the right of the solid line represents the percentage of bootstrapped models considered
 243 to have model violations.

244 This visual summary provides an intuitive way to assess the model's fit and potential
 245 violations, allowing users to quickly grasp the results of the automated analysis.

246 4.5. Modularized infrastructure

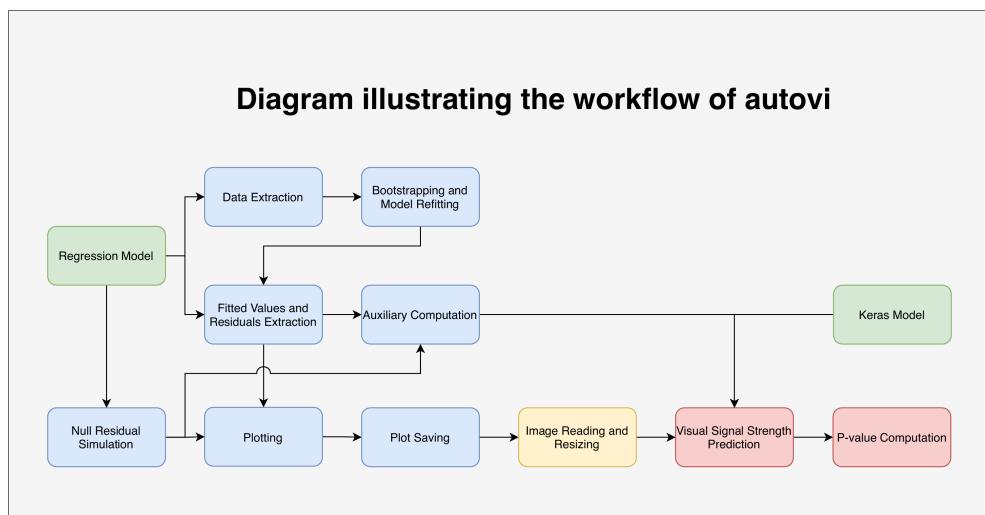


Figure 4. Diagram illustrating the infrastructure of the `extsf{R}` package `autovi`. The modules in green are primary inputs provided by users. Modules in blue are overridable methods that can be modified to accommodate users' specific needs. The module in yellow is a pre-defined non-overridable method. The modules in red are primary outputs of the package.

247 The initial motivation for developing `autovi` was to create a convenient interface for
 248 sharing the models described and trained in [Li et al. \(2024b\)](#). However, recognizing

249 that the classical normal linear regression model represents a restricted class of
250 models, we sought to avoid limiting the potential for future extensions, whether by
251 the original developers or other developers. As a result, the package was designed to
252 function seamlessly with linear regression models with minimal modification and few
253 required arguments, while also accommodating other classes of models through partial
254 infrastructure substitution. This modular and customizable design allows `autovi` to
255 handle a wide range of residual diagnostics tasks.

256 The infrastructure of `autovi` consists of ten core modules: data extraction,
257 bootstrapping and model refitting, fitted values and residuals extraction, auxiliary
258 computation, null residual simulation, plotting, plot saving, image reading and resizing,
259 VSS prediction, and p -value computation. Each module is designed with minimal
260 dependency on the preceding modules, allowing users to customize parts of the
261 infrastructure without affecting its overall integrity. An overview of this infrastructure
262 is illustrated in Figure 4.

263 The package takes regression models and a `Keras` model as primary inputs. Modules for
264 VSS prediction and p -value computation are fixed but accessible via function arguments,
265 using `TensorFlow` for inference and statistical testing. The image loading module is
266 also fixed, using `PIL` to read and resize images based on the `Keras` model's input shape.
267 The remaining seven modules are overridable, allowing users to adapt the workflow as
268 needed. The data extraction module extracts a `data.frame` containing variables used
269 in the regression model. The bootstrapping and refitting module resamples the data
270 and refits the model. The fitted values and residuals extraction module returns these
271 values as a `data.frame`. The auxiliary computation module calculates scagnostics
272 such as monotonicity. The plotting module generates a `ggplot` in a standard format,
273 and the plot saving module exports it at the same resolution as the training images.
274 These modules are described in detail in the package documentation.

275 4.6. Extension to Other Model Classes

276 The `autovi` R package can be extended to accommodate other classes of models
277 beyond linear regression, such as generalized linear models (`glm`). This is achieved
278 by substituting the relevant overridable modules, and if needed, supplying a different
279 `Keras` model.

280 We provide an example of defining a new checker class tailored for Poisson regression
281 using the `glm` framework:

- 282 1. Define a new class using `new_class()` with `AUTO_VI` as the parent class.
- 283 2. Override the necessary methods using `register_method()`. In this example, we
284 use Pearson residuals. To simulate null residuals, we assume the fitted model
285 is correct and the estimated coefficients are accurate. New response values are
286 generated accordingly, and a new model is fitted to this simulated response. Null
287 residuals are then extracted from this refitted model.
- 288 3. Create an alias for the `instantiate()` method of the new class.

```
AUTO_POIS_VI <- new_class(AUTO_VI, class_name = "AUTO_POIS_VI")
register_method(
  AUTO_POIS_VI,
  get_fitted_and_resid = function(fitted_model = self$fitted_model) {
    tibble(.fitted = fitted(fitted_model),
           .resid = resid(fitted_model, type = "pearson"))
  },
  null_method = function(fitted_model = self$fitted_model) {
    dat <- model.frame(fitted_model)
    dat[[1]] <- rpois(nrow(dat), lambda = fitted(fitted_model))
    new_mod <- update(fitted_model, data = dat)
    return(self$get_fitted_and_resid(new_mod))
  }
)
auto_pois_vi <- AUTO_POIS_VI$instantiate
```

- 289 The resulting checker class can be employed analogously to the linear model case
290 described in Section 4.4. For illustration, we fit a Poisson model in which the quadratic
291 term of the predictor x is intentionally omitted. This misspecification manifests as
292 a pronounced U-shaped pattern in the lineup display (see Figure 5), which is also
293 successfully identified by the computer vision model, yielding a p-value substantially
294 below the conventional threshold of 0.05.

```
x <- rnorm(300, sd = 0.5)
y <- rpois(300, lambda = exp(1 + x + x^2))
pois_checker <- auto_pois_vi(
  glm(y ~ x, family = "poisson"),
  keras_model = get_keras_model("vss_phn_32"))
)
pois_checker$plot_lineup()
```

The true residual plot is at position 4.

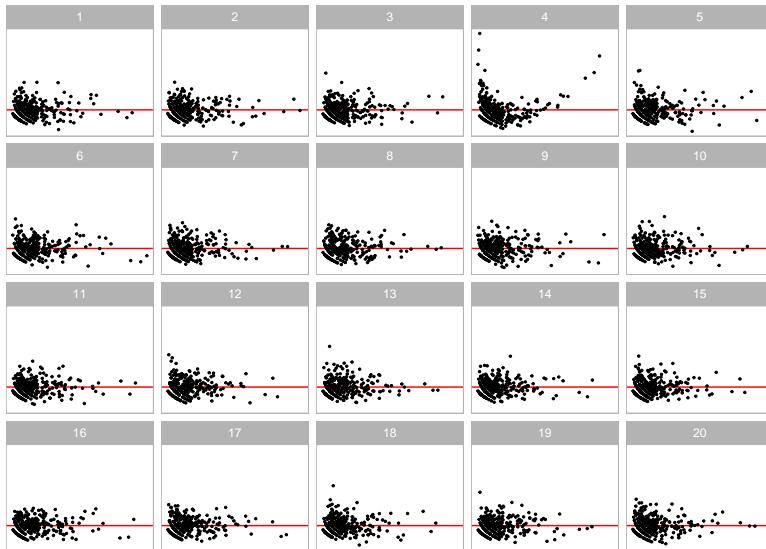


Figure 5. A lineup of residual plots from Poisson generalized linear models, with the true residual plot at position 4, which displays a distinct U-shaped pattern. In contrast, the null plots show characteristics broadly consistent with well-behaved residuals from linear regression models.

```
pois_checker$check()
```

```

295 <AUTO_POIS_VI object>
296 Status:
297   - Fitted model: glm, lm
298   - Keras model: (None, 32, 32, 3) + (None, 5) -> (None, 1)
299     - Output node index: 1
300   - Result:
301     - Observed visual signal strength: 4.875 (p-value = 0.009901)
302     - Null visual signal strength: [100 draws]
303       - Mean: 1.331
304       - Quantiles:
305
306         25%    50%    75%    80%    90%    95%    99%
307         1.035  1.233  1.488  1.644  1.941  2.276  2.639
308
309     - Bootstrapped visual signal strength: [100 draws]
310       - Mean: 5.51 (p-value = 0.009901)

```

311 - Quantiles:

312

25%	50%	75%	80%	90%	95%	99%
5.330	5.505	5.698	5.735	5.830	5.903	6.013

313

314

315

316 - Likelihood ratio: 0.05096 (boot) / 0 (null) = Extremely large

317 It is important to note, however, that the pre-trained computer vision model included
 318 in `autovi`, such as `vss_phn_32` (see `list_keras_models()` for the full list of available
 319 models), was developed specifically for diagnostics of linear regression. Its applicability
 320 to other model classes relies on the assumption that the null residual plots exhibit
 321 characteristics broadly consistent with those of well-behaved linear regression residuals,
 322 that is, residuals should be approximately randomly scattered around zero, display
 323 roughly constant variance across the range of fitted values, and exhibit no discernible
 324 structure or curvature. If these conditions are not met, or if model violations do not
 325 give rise to visually detectable patterns, the validity of the automated diagnostics may
 326 be compromised. In such cases, users are encouraged to train and apply their own
 327 `Keras` models. Detailed guidance on model training and discussion on extending the
 328 methodology to other model classes can be found in [Li et al. \(2024b\)](#).

329 5. Web interface: `autovi.web`

330 The `autovi.web` Shiny application extends the functionality of `autovi` by offering a
 331 user-friendly web interface for automated residual plot assessment. This eliminates the
 332 common challenges associated with software installation, so users can avoid managing
 333 `Python` environments or handling version requirements for `R` libraries. The platform
 334 is cross-platform and accessible on various devices and operating systems, making it
 335 suitable even for users without `R` programming experience. Additionally, updates are
 336 managed centrally, ensuring that users always have access to the latest features. This
 337 section discusses the implementation based on `autovi.web` version 0.1.0.

338 5.1. Implementation

339 The interface `autovi.web` is built using the `shiny` ([Chang et al. 2022](#)) and
 340 `shinydashboard` ([Chang & Borges Ribeiro 2021](#)) `R` packages. Hosted on the
 341 `shinyapps.io` domain, the application is accessible through any modern web browser.
 342 The `R` packages `htmltools` ([Cheng et al. 2024](#)) and `shinycssloaders` ([Sali & Attali](#)

343 2020) are used to render markdown documentation in Shiny application, and for loading
344 animations for Shiny widgets, respectively.

345 Determining the best way to implement the backend was difficult. In our initial
346 planning for `autovi.web`, we considered implementing the entire web application using
347 the `webr` framework (Moon 2020), which would have allowed the entire application to
348 run directly in the user’s browser. However, `webr` does not support packages which use
349 compiled fortran code, which is required by `splancs` (Rowlingson & Diggle 2023), a
350 dependency of `autovi`. In the future, it is possible that a working Emscripten (Zakai
351 2011) version of this package may allow full `webr` support.

352 We also explored the possibility of implementing the web interface using frameworks
353 built on other languages, such as Python. However, server hosting domains that
354 natively support Python servers typically do not have the latest version of R installed.
355 Additionally, calling R from Python is typically done using the `rpy2` Python library
356 (Gautier 2024), but this approach can be awkward when dealing with language syntax
357 related to non-standard evaluation. Another option we considered was renting a server
358 where we could have full control, such as those provided by cloud platforms like
359 Google Cloud Platform (GCP) or Amazon Web Services (AWS). However, deploying
360 and maintaining the server securely requires some expertise. Ultimately, the most
361 practical solution was to use the `shiny` and `shinydashboard` frameworks, which are
362 well-established in the R community and offer a solid foundation for web application
363 development.

364 The server-side configuration of `autovi.web` is carefully designed to support its
365 functionality. Most required Python libraries, including `pillow` and `numpy`, are pre-
366 installed on the server. These libraries are integrated into the Shiny application using
367 the `reticulate` package, which provides an interface between R and Python.

368 Due to shinyapps.io’s resource policy, inactive servers enter sleep mode, clearing the
369 local Python environment. When reactivated for a new session, libraries must be
370 reinstalled. While this ensures a clean environment for each session, it may lead to
371 slightly longer loading times for the first user after a period of inactivity.

372 In contrast to `autovi`, `autovi.web` leverages `TensorFlow.js`, a JavaScript library
373 that allows the execution of machine learning models directly in the browser. This
374 choice enables native browser execution, enhancing compatibility across different user
375 environments, and shifts the computational load from the server to the client-side.

376 TensorFlow.js also offers better scalability and performance, especially when dealing
377 with resource-intensive computer vision models on the web.

378 While autovi requires downloading the pre-trained computer vision models from
379 GitHub, these models in “keras” file format are incompatible with TensorFlow.js.
380 Therefore, we extract and store the model weights in JSON files and include
381 them as extra resources in the Shiny application. When the application initializes,
382 TensorFlow.js rebuilds the computer vision model using these pre-stored weights.

383 To allow communication between TensorFlow.js and other components of the Shiny
384 application, the shinyjs R package ([Attali 2021](#)) is used. This package allows calling
385 custom JavaScript code within the Shiny framework. The specialized JavaScript code for
386 initializing TensorFlow.js and calling TensorFlow.js for VSS prediction is deployed
387 alongside the Shiny application as additional resources.

388 5.2. Usage

389 The workflow of autovi.web is designed to be straightforward, with numbered
390 steps displayed in each panel. There are two example datasets provided by the
391 web application. The single residual plot example uses the dino dataset from the
392 R package **datasauRus** ([Davies, Locke & D'Agostino McGowan 2022](#)). The lineup
393 example uses residuals from a simulated regression model that has a non-linearity
394 issue. We walk through the lineup example to further demonstrate the workflow of
395 the web application.

396 5.2.1. Reading data and setting parameters

397 The user can select to upload data as either a single set of residuals and fitted values
398 in a two (or more) column CSV file or a pre-computed lineup of residuals and null
399 datasets in a three (or more) column CSV file (i.e. multiple sets of residuals and fitted
400 values with a column indicating the set label). Here we illustrate use with lineup
401 example data sets (Figure 6). To use the lineup example data, click the “Use Lineup
402 Example” button. The data status will then update to show the number of rows and
403 columns in the dataset, and the CSV type will automatically be selected to the correct
404 option. Since the example dataset follows the variable naming conventions assumed
405 by the web application, the columns for fitted values, residuals, and labels of residual
406 plots are automatically mapped such that the column named as `.fitted` is mapped
407 to fitted values, `.resid` is mapped to residuals and if applicable, `.sample` to labels of
408 the residual set (middle image). If the user is working with a custom dataset, these

409 options must be set accordingly. Whenever a data containing a lineup, the user must
 410 manually select the label for the true residual plot, otherwise the web application does
 411 not provide all the results. The last step is to click the play button (right image) to
 412 start the assessment.

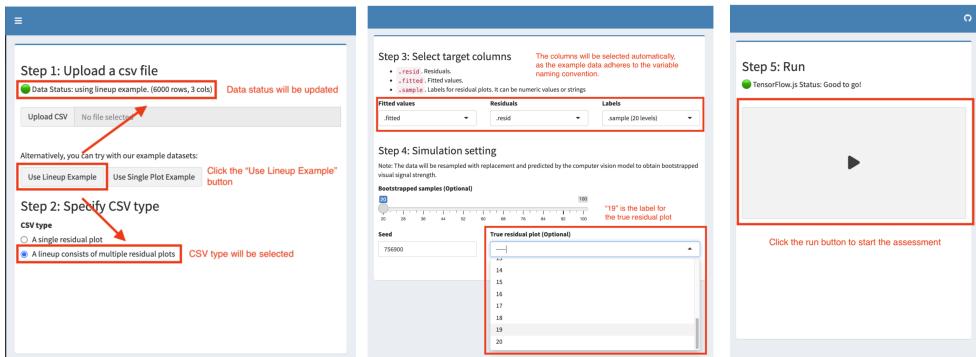


Figure 6. To begin the workflow for `autoovi` using the lineup example dataset, the user clicks the “Use Lineup Example” button (left) to load the example dataset, during which the data status and CSV type will be automatically updated. The user must manually select the label for the true residual plot (middle) to compute further results. The user initiates the assessment of the lineup example data by clicking the run button (right).

413 5.2.2. Results provided

414 Results are provided in multiple panels. The first row of the table Figure 7 is the most
 415 crucial to check, as it provides the VSS and the rank of the true residual plot among
 416 the other plots. The summary text beneath the table provides the *p*-value, which can
 417 be used for quick decision-making. The lineup is for manual inspection, and the user
 418 should see if the true residual plot is visually distinguishable from the other plots, to
 419 confirm if the model violation is serious.

420 The density plot in Figure 8 offers a more robust result, allowing the user to compare
 421 the distribution of bootstrapped VSS with the distribution of null VSS. Finally, the
 422 grayscale attention map (right image) can be used to check if the target visual features,
 423 like the non-linearity present in the lineup example, are captured by the computer
 424 vision model, ensuring the quality of the assessment. The attention map is the gradient
 425 of the model output with respect to the grayscale image input, indicating the sensitivity
 426 of the output to each pixel.

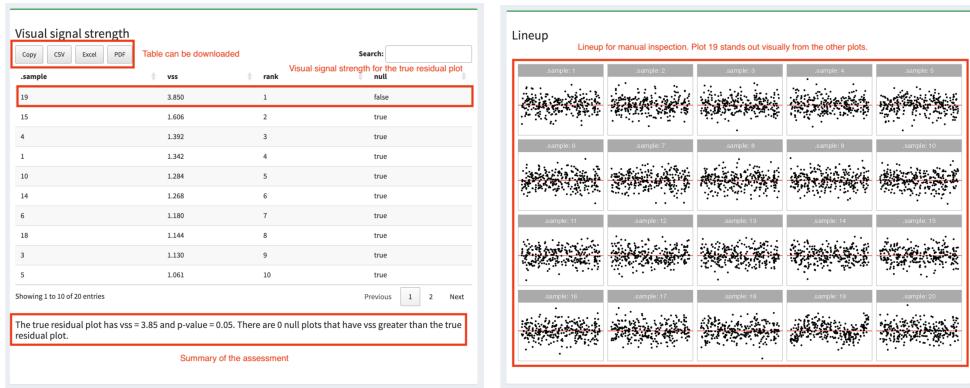


Figure 7. Results for the lineup. The VSS of the true residual plot is displayed in the first row of the table of VSS values for all the null plots (left image), with a summary text beneath the table providing the p -value to aid in decision-making. A lineup of residual plots allows for manual inspection (right image).

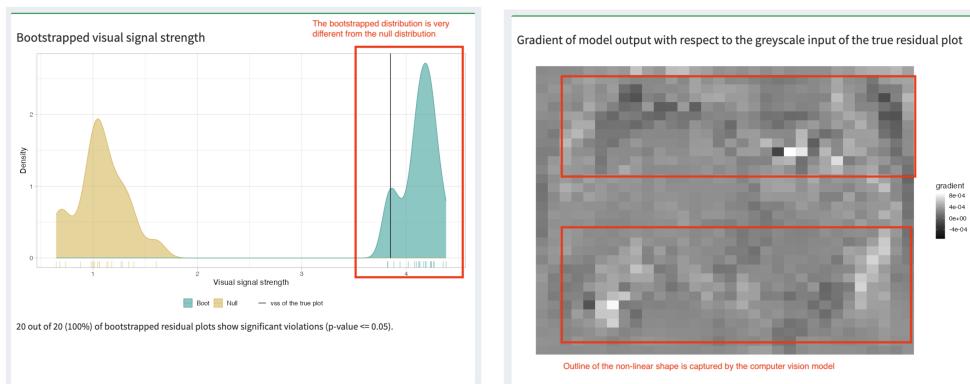


Figure 8. Summaries assessing the strength of the pattern and which elements of the plot contribute. The density plot helps verify if the bootstrapped distribution differs from the null distribution (left image). The attention map (right image) offers insights into whether the computer vision model has captured the intended visual features of the true residual plot.

427

6. Conclusions

428 This paper presents new regression diagnostics software, the R package **autovi** and
 429 its accompanying web interface, **autovi.web**. It addresses a critical gap in the current
 430 landscape of statistical software. While regression tools are widely available, effective
 431 and efficient diagnostic methods have lagged behind, particularly in the field of residual
 432 plot interpretation.

433 The **autovi** R package, introduced in this paper, automates the assessment of residual
 434 plots by incorporating a computer vision model, reducing reliance on time-consuming

435 and potentially inconsistent human interpretation. This automation improves the
436 efficiency of the diagnostic process and promotes consistency in model evaluation
437 across different users and studies.

438 The development of the accompanying Shiny app, `autovi.web`, expands access to these
439 advanced diagnostic tools, by providing a user-friendly interface. It makes automated
440 residual plot assessment accessible to a broader audience, including those who may not
441 have extensive programming experience. This web-based solution effectively addresses
442 the potential barriers to adoption, such as complex dependencies and installation
443 requirements, that are often associated with advanced statistical software.

444 The combination of `autovi` and `autovi.web` offers a comprehensive solution to the
445 challenges of residual plot interpretation in regression analysis. These tools have the
446 potential to significantly improve the quality and consistency of model diagnostics
447 across various fields, from academic research to industry applications. By automating
448 a critical aspect of model evaluation, they allow researchers and analysts to focus more
449 on interpreting results and refining models, rather than grappling with the intricacies
450 of plot assessment.

451 The framework established by `autovi` and `autovi.web` opens up exciting possibilities
452 for further research and development. Future work could explore the extension of these
453 automated assessment techniques to other types of diagnostic plots and statistical
454 models, potentially revolutionizing how we approach statistical inference using visual
455 displays more broadly.

456 7. Resources and supplementary material

457 The current version of `autovi` can be installed from CRAN, and source
458 code for both packages are available at github.com/TengMCing/autovi and
459 github.com/TengMCing/autovi_web respectively. The web interface is available from
460 autoviweb.netlify.app.

461 This paper is reproducibly written using Quarto ([Allaire et al. 2024](#)) powered by
462 Pandoc ([MacFarlane, Krewinkel & Rosenthal 2024](#)) and pdfTeX. The full source code
463 to reproduce this paper is available at github.com/TengMCing/autovi_paper.

464 These R packages were used for the work: tidyverse ([Wickham et al. 2019](#)), lmtest
465 ([Zeileis & Hothorn 2002](#)), kableExtra ([Zhu 2021](#)), patchwork ([Pedersen 2022](#)), rcartocolor
466 ([Nowosad 2018](#)), glue ([Hester & Bryan 2022](#)), here ([Müller 2020](#)), magick ([Ooms 2023](#)),

467 `yardstick` (Kuhn, Vaughan & Hvitfeldt 2024) and `reticulate` (Ushey, Allaire & Tang
 468 2024).

469 **References**

- 470 ABADI, M., AGARWAL, A., BARHAM, P., BREVDO, E., CHEN, Z., CITRO, C., CORRADO, G.S.,
 471 DAVIS, A., DEAN, J., DEVIN, M. et al. (2016). Tensorflow: Large-scale machine learning on
 472 heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* .
- 473 ALLAIRE, J., TEAGUE, C., SCHEIDECKER, C., XIE, Y. & DERVIEUX, C. (2024). Quarto. doi:
 474 10.5281/zenodo.5960048. URL <https://github.com/quarto-dev/quarto-cli>.
- 475 ATTALI, D. (2021). *shinyjs: Easily Improve the User Experience of Your Shiny Apps in Seconds*.
 476 URL <https://CRAN.R-project.org/package=shinyjs>. R package version 2.1.0.
- 477 BALAMUTA, J.J. (2024). *surreal: Create Datasets with Hidden Images in Residual Plots*. URL
 478 <https://CRAN.R-project.org/package=surreal>. R package version 0.0.1.
- 479 BREUSCH, T.S. & PAGAN, A.R. (1979). A simple test for heteroscedasticity and random coefficient
 480 variation. *Econometrica: Journal of the Econometric Society* , 1287–1294.
- 481 BUJA, A., COOK, D., HOFMANN, H., LAWRENCE, M., LEE, E.K., SWAYNE, D.F. & WICKHAM, H.
 482 (2009). Statistical inference for exploratory data analysis and model diagnostics. *Philosophical
 483 Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* **367**,
 484 4361–4383.
- 485 CHANG, W. & BORGES RIBEIRO, B. (2021). *shinydashboard: Create Dashboards with 'Shiny'*.
 486 URL <https://CRAN.R-project.org/package=shinydashboard>. R package version 0.7.2.
- 487 CHANG, W., CHENG, J., ALLAIRE, J., SIEVERT, C., SCHLOERKE, B., XIE, Y., ALLEN, J.,
 488 MCPHERSON, J., DIPERT, A. & BORGES, B. (2022). *shiny: Web Application Framework for
 489 R*. URL <https://CRAN.R-project.org/package=shiny>. R package version 1.7.3.
- 490 CHENG, J., SIEVERT, C., SCHLOERKE, B., CHANG, W., XIE, Y. & ALLEN, J. (2024). *htmltools:
 491 Tools for HTML*. URL <https://CRAN.R-project.org/package=htmltools>. R package version
 492 0.5.8.
- 493 CLARK, A. et al. (2015). Pillow (pil fork) documentation. *readthedocs* .
- 494 COOK, R.D. & WEISBERG, S. (1982). *Residuals and influence in regression*. New York: Chapman
 495 and Hall.
- 496 DAVIES, R., LOCKE, S. & D'AGOSTINO McGOWAN, L. (2022). *datasauRus: Datasets from the
 497 Datasaurus Dozen*. URL <https://CRAN.R-project.org/package=datasauRus>. R package
 498 version 0.1.6.
- 499 GAUTIER, L. (2024). *Python interface to the R language (embedded R)*. URL [https://pypi.org/project/rpy2/](https://pypi.org/

 500 project/rpy2/). Version 3.5.16.
- 501 GOODE, K. & REY, K. (2019). *ggResidpanel: Panels and Interactive Versions of Diagnostic Plots
 502 using 'ggplot2'*. URL <https://CRAN.R-project.org/package=ggResidpanel>. R package version
 503 0.3.0.
- 504 HARTIG, F. (2022). *DHARMa: Residual Diagnostics for Hierarchical (Multi-Level / Mixed)
 505 Regression Models*. URL <https://CRAN.R-project.org/package=DHARMa>. R package
 506 version 0.4.6.
- 507 HEBBALI, A. (2024). *olsrr: Tools for Building OLS Regression Models*. URL [https://CRAN.R-project.org/package=olsrr](https://CRAN.R-

 508 project.org/package=olsrr). R package version 0.6.0.
- 509 HESTER, J. & BRYAN, J. (2022). *glue: Interpreted String Literals*. URL [https://CRAN.R-project.org/package=glue](https://CRAN.R-

 510 project.org/package=glue). R package version 1.6.2.
- 511 JOHNSON, P.E. (2022). *rockchalk: Regression Estimation and Presentation*. URL [https://CRAN.R-project.org/package=rockchalk](https://CRAN.R-

 512 project.org/package=rockchalk). R package version 1.8.157.

- 513 KUHN, M., VAUGHAN, D. & HVITFELDT, E. (2024). *yardstick: Tidy Characterizations of Model*
 514 *Performance*. URL <https://CRAN.R-project.org/package=yardstick>. R package version 1.3.1.
- 515 LI, W. (2024). *bandicoot: Light-weight python-like object-oriented system*. URL <https://CRAN.R->
 516 [project.org/package=bandicoot](https://CRAN.R-project.org/package=bandicoot).
- 517 LI, W., COOK, D., TANAKA, E. & VANDERPLAS, S. (2024a). A plot is worth a thousand tests:
 518 Assessing residual diagnostics with the lineup protocol. *Journal of Computational and*
 519 *Graphical Statistics* **33**, 1497–1511. doi:10.1080/10618600.2024.2344612.
- 520 LI, W., COOK, D., TANAKA, E., VANDERPLAS, S. & ACKERMANN, K. (2024b). Automated
 521 assessment of residual plots with computer vision models. *arXiv preprint arXiv:2411.01001* .
- 522 LONG, J.A. (2022). *jtools: Analysis and Presentation of Social Scientific Data*. URL <https://cran.r->
 523 [project.org/package=jtools](https://cran.r-project.org/package=jtools). R package version 2.2.0.
- 524 LOY, A. & HOFMANN, H. (2014). *Hlmdiag: A suite of diagnostics for hierarchical linear models in*
 525 *r*. *Journal of Statistical Software* **56**, 1–28.
- 526 MACFARLANE, J., KREWINKEL, A. & ROSENTHAL, J. (2024). *Pandoc*. URL <https://github.com/>
 527 [jgm/pandoc](https://github.com/jgm/pandoc).
- 528 MASON, H., LEE, S., LAA, U. & COOK, D. (2022). *cassowaryr: Compute Scagnostics on Pairs of*
 529 *Numeric Variables in a Data Set*. URL <https://CRAN.R-project.org/package=cassowary>. R
 530 package version 2.0.0.
- 531 MOON, K.W. (2020). *webr: Data and Functions for Web-Based Analysis*. URL <https://CRAN.R->
 532 [project.org/package=webr](https://CRAN.R-project.org/package=webr). R package version 0.1.5.
- 533 MÜLLER, K. (2020). *here: A simpler way to find your files*. URL <https://CRAN.R-project.org/>
 534 [package=here](https://CRAN.R-project.org/package=here). R package version 1.0.1.
- 535 NOWOSAD, J. (2018). 'CARTOCOLORs' palettes. URL <https://nowosad.github.io/rkartocolor>. R
 536 package version 1.0.
- 537 OOMS, J. (2023). *magick: Advanced Graphics and Image-Processing in R*. URL <https://CRAN.R->
 538 [project.org/package=magick](https://CRAN.R-project.org/package=magick). R package version 2.7.4.
- 539 PEDERSEN, T.L. (2022). *patchwork: The composer of plots*. URL <https://CRAN.R-project.org/>
 540 [package=patchwork](https://CRAN.R-project.org/package=patchwork). R package version 1.1.2.
- 541 R CORE TEAM (2022). *R: A Language and Environment for Statistical Computing*. R Foundation
 542 for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 543 RAMSEY, J.B. (1969). Tests for specification errors in classical linear least-squares regression
 544 analysis. *Journal of the Royal Statistical Society: Series B (Methodological)* **31**, 350–371.
- 545 REINHART, A. (2024). *regressinator: Simulate and Diagnose (Generalized) Linear Models*. URL
 546 <https://CRAN.R-project.org/package=regressinator>. R package version 0.2.0.
- 547 ROWLINGSON, B. & DIGGLE, P. (2023). *splancs: Spatial and Space-Time Point Pattern Analysis*.
 548 URL <https://CRAN.R-project.org/package=splancs>. R package version 2.01-44.
- 549 SALI, A. & ATTALI, D. (2020). *shinyccsloaders: Add Loading Animations to a 'shiny' Output*
 550 *While It's Recalculating*. URL <https://CRAN.R-project.org/package=shinyccsloaders>. R
 551 package version 1.0.0.
- 552 SHAPIRO, S.S. & WILK, M.B. (1965). An analysis of variance test for normality (complete samples).
 553 *Biometrika* **52**, 591–611.
- 554 USHEY, K., ALLAIRE, J. & TANG, Y. (2024). *reticulate: Interface to 'Python'*. URL <https://CRAN.R->
 555 [project.org/package=reticulate](https://CRAN.R-project.org/package=reticulate). R package version 1.35.0.
- 556 WARTON, D.I. (2023). Global simulation envelopes for diagnostic plots in regression models. *The*
 557 *American Statistician* **77**, 425–431.
- 558 WICKHAM, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
 559 URL <https://ggplot2.tidyverse.org>.
- 560 WICKHAM, H., AVERICK, M., BRYAN, J., CHANG, W., McGOWAN, L.D., FRANÇOIS, R.,
 561 GROLEMUND, G., HAYES, A., HENRY, L., HESTER, J., KUHN, M., PEDERSEN, T.L., MILLER,

- 562 E., BACHE, S.M., MÜLLER, K., OOMS, J., ROBINSON, D., SEIDEL, D.P., SPINU, V.,
563 TAKAHASHI, K., VAUGHAN, D., WILKE, C., WOO, K. & YUTANI, H. (2019). Welcome to
564 the tidyverse. *Journal of Open Source Software* **4**, 1686. doi:10.21105/joss.01686.
- 565 WICKHAM, H., CHOWDHURY, N.R., COOK, D. & HOFMANN, H. (2020). *nullabor: Tools for*
566 *Graphical Inference*. URL <https://CRAN.R-project.org/package=nullabor>. R package version
567 0.3.9.
- 568 ZAKAI, A. (2011). Emscripten: an llvm-to-javascript compiler. In *Proceedings of the ACM*
569 *international conference companion on Object oriented programming systems languages and*
570 *applications companion*. pp. 301–312.
- 571 ZEILEIS, A. & HOTHORN, T. (2002). Diagnostic checking in regression relationships. *R News* **2**,
572 7–10.
- 573 ZHU, H. (2021). *kableExtra: Construct complex table with kable and pipe syntax*. URL
574 <https://CRAN.R-project.org/package=kableExtra>. R package version 1.3.4.