

1 **Automated Residual Plot Assessment with the R Package**
2 **autovi and the Shiny App autovi.web**

3 Weihao Li¹, 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 the adequacy of their 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 `lmtest` (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,
15 and residuals vs. leverage plots. Additional packages like `jtools` (Long 2022), `olsrr`

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: `weihao.li@monash.edu`

(Hebbali 2024), `rockchalk` (Johnson 2022), and `ggResidpanel` (Goode & Rey 2019) deliver similar 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 4 provides a detailed documentation of the `autovi` package, including its usage and infrastructure. Section 5 focuses on the

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

52 2. Definition and computation of visual signal strength

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

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

68 where D_{KL} is given by:

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

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

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

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

78 The computer vision model approximates this mapping from a set of residuals to its
79 corresponding distance measure. It is trained on a large number of synthetic regression
80 models, where the data-generating process is known, allowing the distance measure
81 to be explicitly calculated. The model takes a residual plot as input and outputs the
82 corresponding distance measure. Additional details are provided in [Li et al. \(2024b\)](#).

83 3. Definition and simulation of null and bootstrapped residuals

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

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

95 Our computer vision model is trained to assign lower VSS to null plots and higher VSS
96 to plots that display distinct patterns. Accordingly, statistical testing is performed
97 by computing the proportion of null plots whose VSS equals or exceeds that of the
98 observed residual plot. This proportion serves as a p-value for a one-sided hypothesis
99 test.

100 **Bootstrapped residuals** are obtained by refitting the model on bootstrap samples,
101 which are generated by sampling individual observations with replacement from the
102 original dataset. The residual plots obtained from these refitted models are evaluated
103 using the same computer vision model. The predicted VSS from the bootstrapped
104 plots provide an empirical estimate of the variation in the VSS of the observed
105 residual plot. By examining the proportion of bootstrapped plots that also exhibit

106 significant violations, we can assess whether the original conclusion is robust to
107 sampling variability.

108 **4. R package: autovi**

109 The main purpose of **autovi** is to provide rejection decisions and p -values for testing
110 the null hypothesis (H_0) that the regression model is correctly specified. The package
111 provides automated interpretation of residual plots using computer vision. The name
112 **autovi** stands for **a**utomated **v**isual **i**nference. This functionality can be accessed
113 through the R package **autovi**, or through a web interface, **autovi.web**, which enables
114 users to perform analyses without installing R, Python, or their associated dependencies
115 locally.

116 **4.1. Motivation**

117 Figure 1 shows three sets of plots of residuals against fitted values. The simulated
118 example in (a) might be interpreted as a heteroscedastic pattern, however the
119 automated reading would predict this to have a visual signal strength (VSS) of
120 1.53, with a corresponding p -value of 0.25. This means it would be interpreted as
121 a good residual plot, that there is nothing in the data to indicate a violation of
122 model assumptions. Skewness in the predictor variables is generating the apparent
123 heteroscedasticity, where the smaller variance in residuals at larger fitted values is
124 due to smaller sample size only. The Breusch-Pagan test (Breusch & Pagan 1979) for
125 heteroscedasticity would also not reject this as good residual plot.

126 The data in (b) is generated by fitting a linear model predicting `mpg` based on `hp`
127 using the `datasets::mtcars`. It is a small data set, and there is a hint of nonlinear
128 structure not captured by the model. The automated plot reading would predict a
129 VSS of 3.57, which has a p -value less than 0.05. That is, the nonlinear structure is
130 most likely real, and indicates a problem with the model. The conventional test, a
131 Ramsey Regression Equation Specification Error Test (RESET) (Ramsey 1969) would
132 also strongly detect the nonlinearity.

133 The third example is generated using the `surreal` package (Balamuta 2024) where
134 structured residuals are hidden in data, to be revealed if the correct model is specified.
135 Here a quote based on Tukey is used as the residual structure “visual summaries focus
136 on unexpected values”. The automated plot reading predicts the VSS to be 5.87, with
137 a p -value less than 0.05. This structure is blindingly obvious visually, but a RESET
138 test for nonlinear structure would not report a problem. (It would be detected by

139 a Breusch-Pagan for heteroscedasticity and also Shapiro-Wilk test ([Shapiro & Wilk](#)
 140 [1965](#)) for non-normality.)

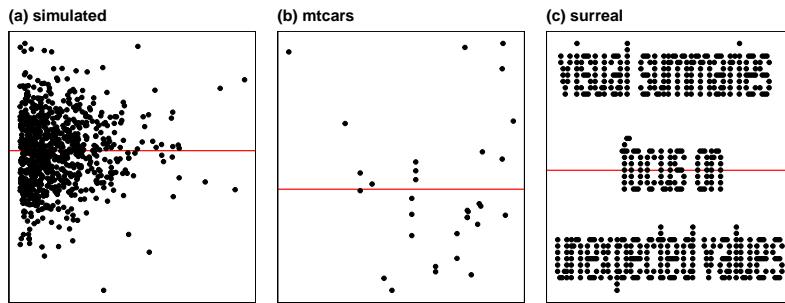


Figure 1. Reading residual plots can be a difficult task, particularly for students new to statistical modeling. The `autovi` package makes it easier. Here are three examples of residual plots, which may appear to have structure. According to `autovi`, the visual signal strengths (VSS) of these three 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.

141 4.2. Implementation

142 The `autovi` package is built on the `bandicoot` object-oriented programming (OOP)
 143 system ([Li 2024](#)), marking a departure from R's traditional S3 generic system. This
 144 OOP architecture enhances flexibility and modularity, allowing users to redefine key
 145 functions through method overriding.

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

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

160 **4.3. Installation**

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

168 We recommend using the Shiny app `autovi.web` if users encounter installation
 169 problems.

170 **4.4. Usage**

171 **4.4.1. Numerical summary**

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

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

179 The code to do this is:

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

180 It produces the following summary:

```
181
182 -- <AUTO_VI object>
183 Status:
184 - Fitted model: lm
185 - Keras model: (None, 32, 32, 3) + (None, 5) -> (None, 1)
186   - Output node index: 1
187 - Result:
```

- Observed visual signal strength: 3.16 (p-value = 0.0396)
- Null visual signal strength: [100 draws]
 - Mean: 1.274
 - Quantiles:

25%	50%	75%	80%	90%	95%	99%
0.802	1.111	1.575	1.666	1.919	2.657	3.348

- Bootstrapped visual signal strength: [100 draws]
 - Mean: 2.795 (p-value = 0.05941)
 - Quantiles:

25%	50%	75%	80%	90%	95%	99%
2.455	2.941	3.177	3.300	3.474	3.537	3.668

- Likelihood ratio: 0.7333 (boot) / 0.06284 (null) = 11.67

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

4.4.2. Visual summary

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

```
checker$plot_pair()
```

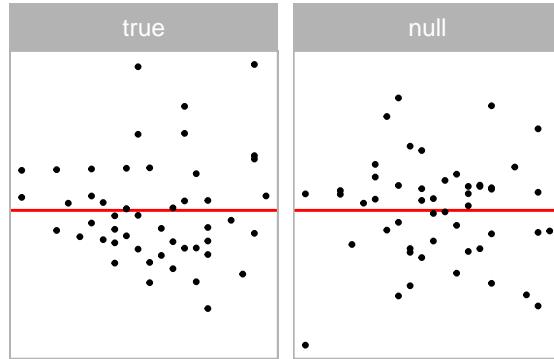


Figure 2. True plot alongside one null plot, for quick comparison.

- 216 The `plot_pair()` method (Figure 2) displays the true residual plot on the left and a
 217 single null plot on the right. If a full lineup was shown, the true residual plot would
 218 be embedded in a page of null plots. Users should look for any distinct visual patterns
 219 in the true residual plot that are absent in the null plot. Running these functions
 220 multiple times can help any visual suspicions, as each execution generates new random
 221 null plots for comparison.
- 222 The package offers a straightforward visualization of the assessment result through
 223 the `summary_plot()` function.

```
checker$summary_plot()
```

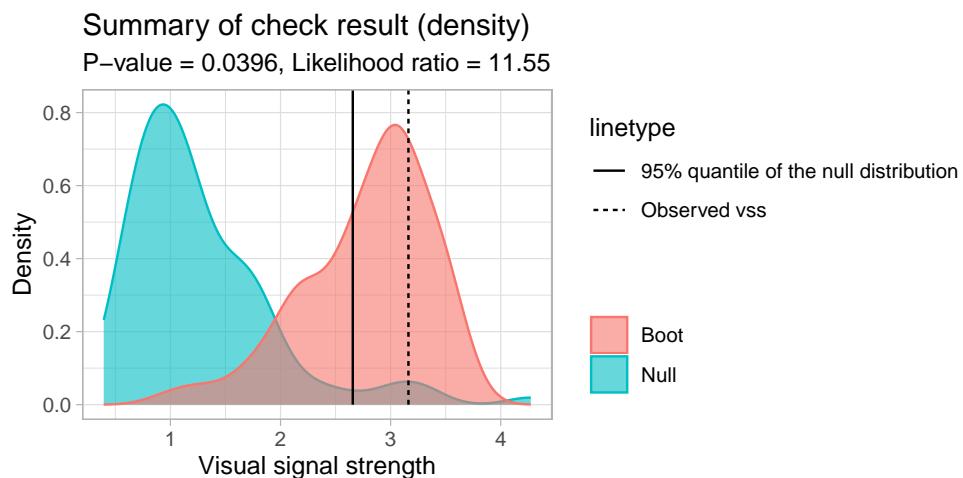


Figure 3. Summary plot comparing the densities of VSS for bootstrapped residual samples (red) relative to VSS for null plots (blue).

- 224 In the result, shown in Figure 3, the blue area represents the density of VSS for null
 225 residual plots, while the red area shows the density for bootstrapped residual plots.
 226 The dashed line indicates the VSS of the true residual plot, and the solid line marks
 227 the critical value at a 95% significance level. The p -value and the likelihood ratio are
 228 displayed in the subtitle. The likelihood ratio represents the ratio of the likelihood
 229 of observing the VSS of the true residual plot from the bootstrapped distribution
 230 compared to the null distribution.
- 231 Interpreting the plot involves several key aspects. If the dashed line falls to the right of
 232 the solid line, it suggests rejecting the null hypothesis. The degree of overlap between
 233 the red and blue areas indicates similarity between the true residual plot and null
 234 plots; greater overlap suggests more similarity. Lastly, the portion of the red area to
 235 the right of the solid line represents the percentage of bootstrapped models considered
 236 to have model violations.
- 237 This visual summary provides an intuitive way to assess the model's fit and potential
 238 violations, allowing users to quickly grasp the results of the automated analysis.

239 **4.5. Modularized infrastructure**

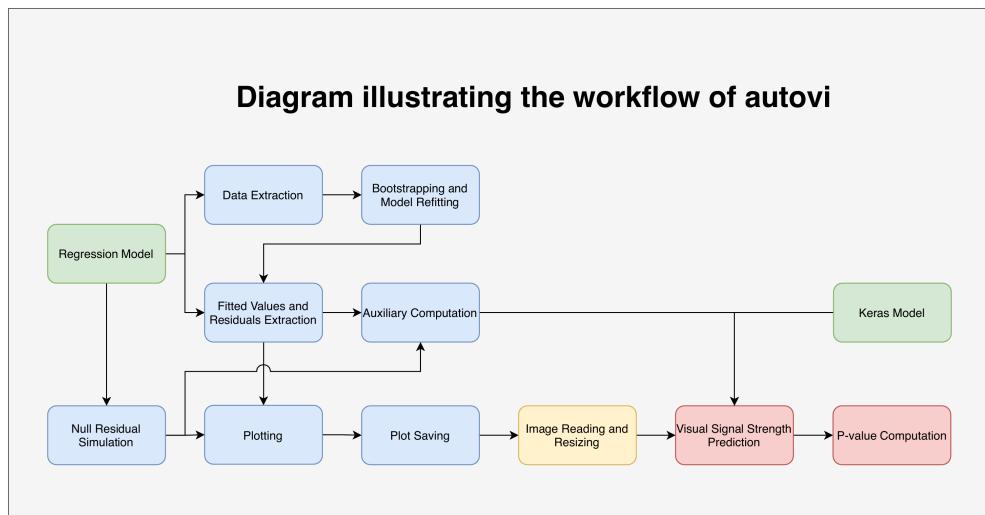


Figure 4. Diagram illustrating the infrastructure of the R package `autovi`. The modules in green are primary inputs provided by users. Modules in blue are overridable methods that can be modified to accommodate users' specific needs. The module in yellow is a pre-defined non-overridable method. The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

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

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

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

256 The modules for VSS prediction and *p*-value computation are predefined and cannot be
257 overridden, although users can interact with them directly through function arguments.
258 Similarly, the image reading and resizing module is fixed but will adapt to different
259 Keras models by checking their input shapes. The remaining seven modules are
260 designed to be overridable, enabling users to tailor the infrastructure to their specific
261 needs. These modules are discussed in detail in the package documentation.

262 4.6. Extension to Other Model Classes

263 The `autovi` R package can be extended to accommodate other classes of models
264 beyond linear regression, such as generalized linear models (`glm`). This is achieved by
265 substituting the relevant overridable modules, as illustrated in Figure 4.

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

- 268 1. Define a new class using `new_class()` with `AUTO_VI` as the parent class.
- 269 2. Override the necessary methods using `register_method()`. In this example, we
270 use Pearson residuals. To simulate null residuals, we assume the fitted model
271 is correct and the estimated coefficients are accurate. New response values are
272 generated accordingly, and a new model is fitted to this simulated response. Null
273 residuals are then extracted from this refitted model.
- 274 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

```

- 275 The resulting checker class can be employed analogously to the linear model case
 276 described in Section 4.4. For illustration, we fit a Poisson model in which the quadratic
 277 term of the predictor x is intentionally omitted. This misspecification manifests as a
 278 pronounced U-shaped pattern in the lineup display, which is also successfully identified
 279 by the computer vision model, yielding a p-value substantially below the conventional
 280 threshold of 0.05.
- 281 It is important to note, however, that the pre-trained computer vision model was
 282 developed specifically for diagnostic evaluation of linear regression. Its applicability
 283 to other model classes relies on the assumption that the null residual plots exhibit
 284 characteristics broadly consistent with those of well-behaved linear regression residuals,
 285 that is, residuals should be approximately randomly scattered around zero, display
 286 roughly constant variance across the range of fitted values, and exhibit no discernible
 287 structure or curvature. If these conditions are not met, or if model violations do not
 288 give rise to visually detectable patterns, the validity of the automated diagnostics may
 289 be compromised. For further discussion on extending the methodology to other model
 290 classes, refer to Li et al. (2024b).

```

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"),

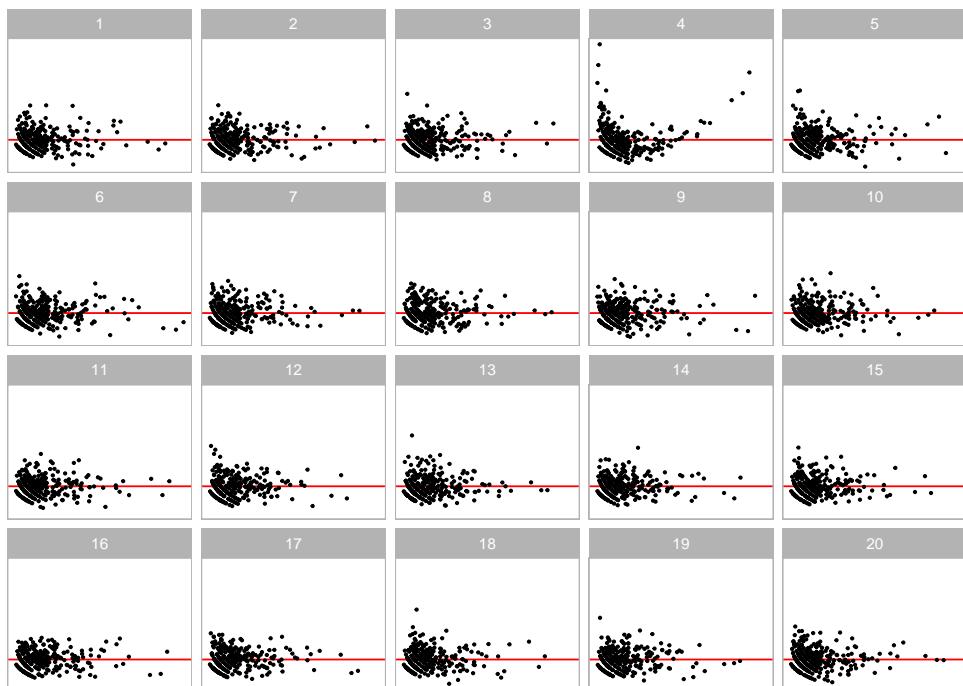
```

```

    keras_model = get_keras_model("vss_phn_32")
)
pois_checker$plot_lineup()

```

The true residual plot is at position 4.



291

```
pois_checker$check()
```

292

```

293 -- <AUTO_POIS_VI object>
294 Status:
295 - Fitted model: glm, lm
296 - Keras model: (None, 32, 32, 3) + (None, 5) -> (None, 1)
297     - Output node index: 1
298 - Result:
299     - Observed visual signal strength: 4.875 (p-value = 0.009901)
300     - Null visual signal strength: [100 draws]
301         - Mean: 1.331
302         - Quantiles:

```



```

304      25%   50%   75%   80%   90%   95%   99%
305      1.035  1.233  1.488  1.644  1.941  2.276  2.639
306
307  - Bootstrapped visual signal strength: [100 draws]
308    - Mean: 5.51 (p-value = 0.009901)
309    - Quantiles:
310
311      25%   50%   75%   80%   90%   95%   99%
312      5.330  5.505  5.698  5.735  5.830  5.903  6.013
313
314  - Likelihood ratio: 0.05096 (boot) / 0 (null) = Extremely large

```

5. Web interface: autovi.web

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

5.1. Implementation

The interface `autovi.web` is built using the `shiny` (Chang et al. 2022) and `shinydashboard` (Chang & Borges Ribeiro 2021) R packages. Hosted on the `shinyapps.io` domain, the application is accessible through any modern web browser. The R packages `htmltools` (Cheng et al. 2024) and `shinycssloaders` (Sali & Attali 2020) are used to render markdown documentation in shiny application, and for loading animations for shiny widgets, respectively.

Determining the best way to implement the backend was difficult. In our initial planning for `autovi.web`, we considered implementing the entire web application using the `webr` framework (Moon 2020), which would have allowed the entire application to run directly in the user's browser. However, `webr` does not support packages which use compiled fortran code, which is required by `splancs` (Rowlingson & Diggle 2023), a

336 dependency of `autovi`. In the future, it is possible that a working Emscripten ([Zakai
337 2011](#)) version of this package may allow full `webr` support.

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

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

354 Due to the resource allocation policy of shinyapps.io, the server enters a sleep mode
355 during periods of inactivity, resulting in the clearing of the local Python virtual
356 environment. Consequently, when the application “wakes up” for a new user session,
357 these libraries need to be reinstalled. While this ensures a clean environment for each
358 session, it may lead to slightly longer loading times for the first user after a period of
359 inactivity.

360 In contrast to `autovi`, `autovi.web` leverages `TensorFlow.js`, a JavaScript library
361 that allows the execution of machine learning models directly in the browser. This
362 choice enables native browser execution, enhancing compatibility across different user
363 environments, and shifts the computational load from the server to the client-side.
364 `TensorFlow.js` also offers better scalability and performance, especially when dealing
365 with resource-intensive computer vision models on the web.

366 While `autovi` requires downloading the pre-trained computer vision models from
367 GitHub, these models in “.keras” file format are incompatible with `TensorFlow.js`.
368 Therefore, we extract and store the model weights in JSON files and include

369 them as extra resources in the Shiny application. When the application initializes,
370 `TensorFlow.js` rebuilds the computer vision model using these pre-stored weights.

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

376 5.2. Usage

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

384 5.2.1. Reading data and setting parameters

385 The user can select to upload data as either a single set of residuals and fitted values
386 in a two (or more) column CSV file or a pre-computed lineup of residuals and null
387 datasets in a three (or more) column CSV file (i.e. multiple sets of residuals and fitted
388 values with a column indicating the set label). Here we illustrate use with lineup
389 example data sets (Figure 5). To use the lineup example data, click the “Use Lineup
390 Example” button. The data status will then update to show the number of rows and
391 columns in the dataset, and the CSV type will automatically be selected to the correct
392 option. Since the example dataset follows the variable naming conventions assumed
393 by the web application, the columns for fitted values, residuals, and labels of residual
394 plots are automatically mapped such that the column named as `.fitted` is mapped
395 to fitted values, `.resid` is mapped to residuals and if applicable, `.sample` to labels of
396 the residual set (middle image). If the user is working with a custom dataset, these
397 options must be set accordingly. Whenever a data containing a lineup, the user must
398 manually select the label for the true residual plot, otherwise the web application does
399 not provide all the results. The last step is to click the play button (right image) to
400 start the assessment.

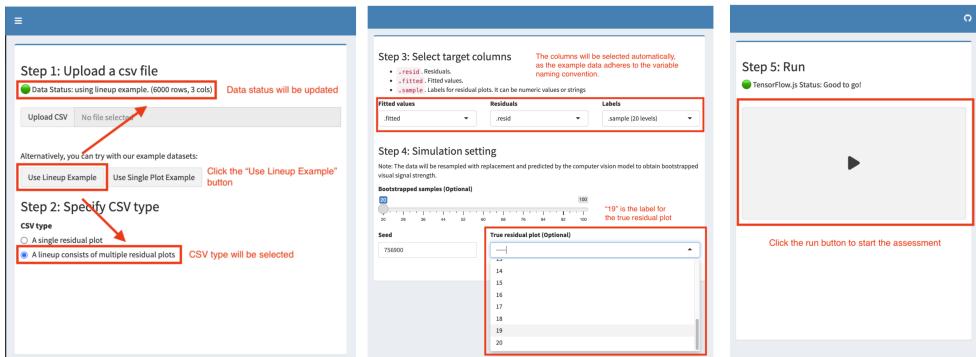


Figure 5. To begin the workflow for `autovi` using the lineup example dataset, the user clicks the “Use Lineup Example” button (left) to load the example dataset, during which the data status and CSV type will be automatically updated. The user must manually select the label for the true residual plot (middle) to compute further results. The user initiates the assessment of the lineup example data by clicking the run button (right).

401 5.2.2. Results provided

402 Results are provided in multiple panels. The first row of the table Figure 6 is the most
 403 crucial to check, as it provides the VSS and the rank of the true residual plot among
 404 the other plots. The summary text beneath the table provides the p -value, which can
 405 be used for quick decision-making. The lineup is for manual inspection, and the user
 406 should see if the true residual plot is visually distinguishable from the other plots, to
 407 confirm if the model violation is serious.

408 The density plot in Figure 7 offers a more robust result, allowing the user to compare
 409 the distribution of bootstrapped VSS with the distribution of null VSS. Finally, the
 410 grayscale attention map (right image) can be used to check if the target visual features,
 411 like the non-linearity present in the lineup example, are captured by the computer
 412 vision model, ensuring the quality of the assessment.

413 6. Conclusions

414 This paper presents new regression diagnostics software, the R package `autovi` and
 415 its accompanying web interface, `autovi.web`. It addresses a critical gap in the current
 416 landscape of statistical software. While regression tools are widely available, effective
 417 and efficient diagnostic methods have lagged behind, particularly in the field of residual
 418 plot interpretation.

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

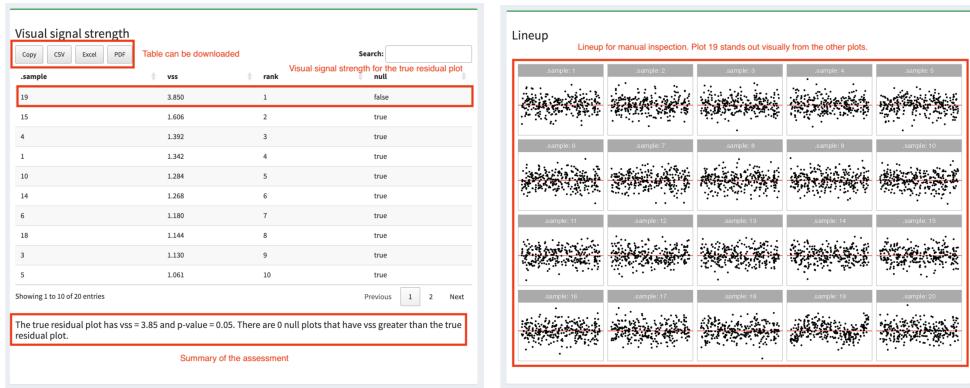


Figure 6. Results for the lineup. The VSS of the true residual plot is displayed in the first row of the table of VSS values for all the null plots (left image), with a summary text beneath the table providing the p -value to aid in decision-making. A lineup of residual plots allows for manual inspection (right image).

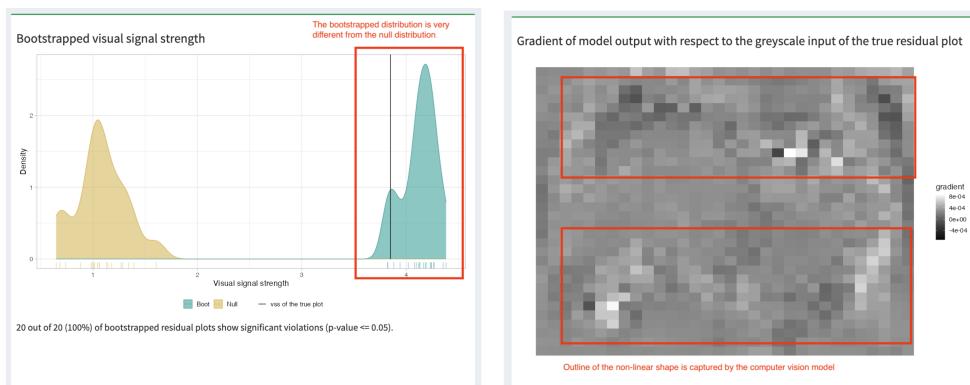


Figure 7. Summaries assessing the strength of the pattern and which elements of the plot contribute. The density plot helps verify if the bootstrapped distribution differs from the null distribution (left image). The attention map (right image) offers insights into whether the computer vision model has captured the intended visual features of the true residual plot.

- 421 and potentially inconsistent human interpretation. This automation improves the
 422 efficiency of the diagnostic process and promotes consistency in model evaluation
 423 across different users and studies.
- 424 The development of the accompanying Shiny app, `autovi.web`, expands access to these
 425 advanced diagnostic tools, by providing a user-friendly interface. It makes automated
 426 residual plot assessment accessible to a broader audience, including those who may not
 427 have extensive programming experience. This web-based solution effectively addresses

428 the potential barriers to adoption, such as complex dependencies and installation
 429 requirements, that are often associated with advanced statistical software.

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

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

442 7. Resources and supplementary material

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

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

450 These R packages were used for the work: `tidyverse` ([Wickham et al. 2019](#)), `lmtest`
 451 ([Zeileis & Hothorn 2002](#)), `kableExtra` ([Zhu 2021](#)), `patchwork` ([Pedersen 2022](#)),
 452 `rcartocolor` ([Nowosad 2018](#)), `glue` ([Hester & Bryan 2022](#)), `here` ([Müller 2020](#)),
 453 `magick` ([Ooms 2023](#)), `yardstick` ([Kuhn, Vaughan & Hvitfeldt 2024](#)) and `reticulate`
 454 ([Ushey, Allaire & Tang 2024](#)).

455 References

- 456 ABADI, M., AGARWAL, A., BARHAM, P., BREVDO, E., CHEN, Z., CITRO, C., CORRADO, G.S.,
 457 DAVIS, A., DEAN, J., DEVIN, M. et al. (2016). Tensorflow: Large-scale machine learning on
 458 heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* .
 459 ALLAIRE, J., TEAGUE, C., SCHEIDEGGER, C., XIE, Y. & DERVIEUX, C. (2024). Quarto. doi:
 460 10.5281/zenodo.5960048. URL <https://github.com/quarto-dev/quarto-cli>.

- 461 ATTALI, D. (2021). *shinyjs: Easily Improve the User Experience of Your Shiny Apps in Seconds*.
 462 URL <https://CRAN.R-project.org/package=shinyjs>. R package version 2.1.0.
- 463 BALAMUTA, J.J. (2024). *surreal: Create Datasets with Hidden Images in Residual Plots*. URL
 464 <https://CRAN.R-project.org/package=surreal>. R package version 0.0.1.
- 465 BREUSCH, T.S. & PAGAN, A.R. (1979). A simple test for heteroscedasticity and random coefficient
 466 variation. *Econometrica: Journal of the Econometric Society* , 1287–1294.
- 467 BUJA, A., COOK, D., HOFMANN, H., LAWRENCE, M., LEE, E.K., SWAYNE, D.F. & WICKHAM, H.
 468 (2009). Statistical inference for exploratory data analysis and model diagnostics. *Philosophical
 469 Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* **367**,
 470 4361–4383.
- 471 CHANG, W. & BORGES RIBEIRO, B. (2021). *shinydashboard: Create Dashboards with 'Shiny'*.
 472 URL <https://CRAN.R-project.org/package=shinydashboard>. R package version 0.7.2.
- 473 CHANG, W., CHENG, J., ALLAIRE, J., SIEVERT, C., SCHLOERKE, B., XIE, Y., ALLEN, J.,
 474 MCPHERSON, J., DIPERT, A. & BORGES, B. (2022). *shiny: Web Application Framework for
 475 R*. URL <https://CRAN.R-project.org/package=shiny>. R package version 1.7.3.
- 476 CHENG, J., SIEVERT, C., SCHLOERKE, B., CHANG, W., XIE, Y. & ALLEN, J. (2024). *htmltools:
 477 Tools for HTML*. URL <https://CRAN.R-project.org/package=htmltools>. R package version
 478 0.5.8.
- 479 CLARK, A. et al. (2015). Pillow (pil fork) documentation. *readthedocs* .
- 480 COOK, R.D. & WEISBERG, S. (1982). *Residuals and influence in regression*. New York: Chapman
 481 and Hall.
- 482 DAVIES, R., LOCKE, S. & D'AGOSTINO McGOWAN, L. (2022). *datasauRus: Datasets from the
 483 Datasaurus Dozen*. URL <https://CRAN.R-project.org/package=datasauRus>. R package
 484 version 0.1.6.
- 485 GAUTIER, L. (2024). *Python interface to the R language (embedded R)*. URL [https://pypi.org/project/rpy2/](https://pypi.org/

 486 project/rpy2/). Version 3.5.16.
- 487 GOODE, K. & REY, K. (2019). *ggResidpanel: Panels and Interactive Versions of Diagnostic Plots
 488 using 'ggplot2'*. URL <https://CRAN.R-project.org/package=ggResidpanel>. R package version
 489 0.3.0.
- 490 HARTIG, F. (2022). *DHARMA: Residual Diagnostics for Hierarchical (Multi-Level / Mixed)
 491 Regression Models*. URL <https://CRAN.R-project.org/package=DHARMA>. R package
 492 version 0.4.6.
- 493 HEBBALI, A. (2024). *olsrr: Tools for Building OLS Regression Models*. URL [https://CRAN.R-project.org/package=olsrr](https://CRAN.R-

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

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

 498 project.org/package=rockchalk). R package version 1.8.157.
- 499 KUHN, M., VAUGHAN, D. & HVITFELDT, E. (2024). *yardstick: Tidy Characterizations of Model
 500 Performance*. URL <https://CRAN.R-project.org/package=yardstick>. R package version 1.3.1.
- 501 LI, W. (2024). *bandicoot: Light-weight python-like object-oriented system*. URL [https://CRAN.R-project.org/package=bandicoot](https://CRAN.R-

 502 project.org/package=bandicoot).
- 503 LI, W., COOK, D., TANAKA, E. & VANDERPLAS, S. (2024a). A plot is worth a thousand tests:
 504 Assessing residual diagnostics with the lineup protocol. *Journal of Computational and
 505 Graphical Statistics* **33**, 1497–1511. doi:10.1080/10618600.2024.2344612.
- 506 LI, W., COOK, D., TANAKA, E., VANDERPLAS, S. & ACKERMANN, K. (2024b). Automated
 507 assessment of residual plots with computer vision models. *arXiv preprint arXiv:2411.01001* .
- 508 LONG, J.A. (2022). *jtools: Analysis and Presentation of Social Scientific Data*. URL [https://cran.r-project.org/package=jtools](https://cran.r-

 509 project.org/package=jtools). R package version 2.2.0.

- 510 LOY, A. & HOFMANN, H. (2014). Hlmdiag: A suite of diagnostics for hierarchical linear models in
 511 r. *Journal of Statistical Software* **56**, 1–28.
- 512 MACFARLANE, J., KREWINKEL, A. & ROSENTHAL, J. (2024). Pandoc. URL <https://github.com/jgm/pandoc>.
- 514 MASON, H., LEE, S., LAA, U. & COOK, D. (2022). cassowaryr: Compute Scagnostics on Pairs of
 515 Numeric Variables in a Data Set. URL <https://CRAN.R-project.org/package=cassowary>. R
 516 package version 2.0.0.
- 517 MOON, K.W. (2020). webr: Data and Functions for Web-Based Analysis. URL <https://CRAN.R-project.org/package=webr>. R package version 0.1.5.
- 519 MÜLLER, K. (2020). here: A simpler way to find your files. URL <https://CRAN.R-project.org/package=here>. R package version 1.0.1.
- 521 NOWOSAD, J. (2018). 'CARTOCOLORs' palettes. URL <https://nowosad.github.io/rkartocolor>. R
 522 package version 1.0.
- 523 OOMS, J. (2023). magick: Advanced Graphics and Image-Processing in R. URL <https://CRAN.R-project.org/package=magick>. R package version 2.7.4.
- 525 PEDERSEN, T.L. (2022). patchwork: The composer of plots. URL <https://CRAN.R-project.org/package=patchwork>. R package version 1.1.2.
- 527 R CORE TEAM (2022). R: A Language and Environment for Statistical Computing. R Foundation
 528 for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 529 RAMSEY, J.B. (1969). Tests for specification errors in classical linear least-squares regression
 530 analysis. *Journal of the Royal Statistical Society: Series B (Methodological)* **31**, 350–371.
- 531 REINHART, A. (2024). regressinator: Simulate and Diagnose (Generalized) Linear Models. URL
 532 <https://CRAN.R-project.org/package=regressinator>. R package version 0.2.0.
- 533 ROWLINGSON, B. & DIGGLE, P. (2023). splancs: Spatial and Space-Time Point Pattern Analysis.
 534 URL <https://CRAN.R-project.org/package=splancs>. R package version 2.01-44.
- 535 SALI, A. & ATTALI, D. (2020). shinycssloaders: Add Loading Animations to a 'shiny' Output
 536 While It's Recalculating. URL <https://CRAN.R-project.org/package=shinycssloaders>. R
 537 package version 1.0.0.
- 538 SHAPIRO, S.S. & WILK, M.B. (1965). An analysis of variance test for normality (complete samples).
 539 *Biometrika* **52**, 591–611.
- 540 USHEY, K., ALLAIRE, J. & TANG, Y. (2024). reticulate: Interface to 'Python'. URL <https://CRAN.R-project.org/package=reticulate>. R package version 1.35.0.
- 542 WARTON, D.I. (2023). Global simulation envelopes for diagnostic plots in regression models. *The
 543 American Statistician* **77**, 425–431.
- 544 WICKHAM, H. (2016). ggplot2: Elegant graphics for data analysis. Springer-Verlag New York.
 545 URL <https://ggplot2.tidyverse.org>.
- 546 WICKHAM, H., AVERICK, M., BRYAN, J., CHANG, W., McGOWAN, L.D., FRANÇOIS, R.,
 547 GROLEMUND, G., HAYES, A., HENRY, L., HESTER, J., KUHN, M., PEDERSEN, T.L., MILLER,
 548 E., BACHE, S.M., MÜLLER, K., OOMS, J., ROBINSON, D., SEIDEL, D.P., SPINU, V.,
 549 TAKAHASHI, K., VAUGHAN, D., WILKE, C., WOO, K. & YUTANI, H. (2019). Welcome to
 550 the tidyverse. *Journal of Open Source Software* **4**, 1686. doi:10.21105/joss.01686.
- 551 WICKHAM, H., CHOWDHURY, N.R., COOK, D. & HOFMANN, H. (2020). nullabor: Tools for
 552 Graphical Inference. URL <https://CRAN.R-project.org/package=nullabor>. R package version
 553 0.3.9.
- 554 ZAKAI, A. (2011). Emscripten: an llvm-to-javascript compiler. In *Proceedings of the ACM
 555 international conference companion on Object oriented programming systems languages and
 556 applications companion*. pp. 301–312.
- 557 ZEILEIS, A. & HOTHORN, T. (2002). Diagnostic checking in regression relationships. *R News* **2**,
 558 7–10.

- 559 ZHU, H. (2021). *kableExtra: Construct complex table with kable and pipe syntax.* URL
560 <https://CRAN.R-project.org/package=kableExtra>. R package version 1.3.4.