

# pcvr: An R package and tutorials for guided statistical analysis of plant phenotyping data

Josh Sumner<sup>a</sup>, Noah Fahlgren<sup>a</sup>, and Katherine M. Murphy<sup>a</sup>

<sup>a</sup>Donald Danforth Plant Science Center, Saint Louis, MO

## ABSTRACT

**PlantCV** is a Python-based image analysis tool that lowers the barrier to entry for complex image analysis workflows in plant phenotyping. To provide support for subsequent analysis steps of measured trait data we have developed **pcvr**, an R package to assist in common plant phenotyping analyses. The goal of **pcvr** is to make common statistical analyses both easier and more consistent and to lower the barrier to entry for useful Bayesian methods. Here we demonstrate three pieces of a possible analysis covering single value trait analysis, longitudinal modeling, and multi-value trait analysis.

**Keywords:** NAPPN 2024, plant phenotyping, R, Statistics

## 1. INTRODUCTION

High-throughput, image-based phenotyping allows scientists to collect large amounts of data quickly, non-destructively, and accurately. **PlantCV**<sup>1</sup> is a Python-based image analysis tool that lowers the barrier to entry for complex image analysis workflows. **PlantCV** returns phenotype data as numeric values which users generally analyze using R.<sup>2</sup> To provide support for the statistical analysis following a **PlantCV** workflow we developed **pcvr**, an R package to use **PlantCV** output data. **pcvr** aims to make answering common analysis questions easier and to lower the barrier to entry for select Bayesian methods. The **pcvr** package is available on [GitHub](#) where several [tutorials](#) highlighting specific topics or features are also available. **pcvr** includes functions to support reading in large **PlantCV** datasets, variance partitioning in single value traits, outlier removal, single value trait statistical comparisons, longitudinal analysis, multi-value trait analysis, a variety of data visualizations, pseudo water use efficiency calculations, and comparisons of relative stress tolerances. The package also includes vignettes that demonstrate potential analysis workflows for a small dataset collected using the Bellwether Phenotyping Facility at the Donald Danforth Plant Science Center (RRID:SCR\_019049). A subset of the Bellwether vignette's content is shown in this paper, covering single value trait analysis at one time point, longitudinal analysis of single value traits, and multi-value trait analysis.

## 2. SINGLE VALUE TRAIT ANALYSIS

The most commonly used **PlantCV** output data are single value traits. These are phenotypes that can be described by a single number for each image such as plant area, height, width, or number of leaves. These phenotypes are generally used either to compare groups at a single time point or to compare groups over time in longitudinal analysis. **pcvr** includes support for both analyses using both frequentist and Bayesian methods.

Image-based phenotyping allows for non-destructive, repeat sampling of many plants, which is beneficial for measuring temporal responses. However, longitudinal data requires additional considerations when doing statistical analysis. When longitudinal data is collected the first question that a researcher poses tends to concern whether a difference is seen at the end of the experiment. In either case, a single time point comparison will be used to test differences between groups.

---

Further author information: (Send correspondence to Josh Sumner)

Josh Sumner: E-mail: [jsumner@danforthcenter.org](mailto:jsumner@danforthcenter.org), ORCID: 0000-0002-3399-9063

The most common single timepoint hypothesis is that means between groups are different. The way that we might express such a hypothesis verbally is often "How likely is it that group A has a mean lower than group B's mean?" Typically Welch's T-tests or Wilcoxon rank sum tests are used to answer these questions. These frequentist difference of means tests are supported through `ggpubr` in `pcvr::pcvBox` as shown in Figure 1. The P value in Figure 1 is below any reasonable cutoff so our results are declared statistically significant, but our hypothesis as originally stated has not been addressed. The P value is the probability that at least the observed difference would be seen if these two samples are actually drawn from the same population. We cannot remove the condition of there not being a true effect from the P value, so we have not actually answered the question as we originally asked it.

To answer our question as asked we turn to Bayesian statistics. In `pcvr` the `conjugate` function performs Bayesian tests using the method of moments and conjugate prior distributions.<sup>3</sup> Several distributions and hypotheses are supported by `conjugate`. In Figure 2 the "t" method is used to compare the distribution of means from Gaussian data with  $N(50, 15)$  priors. The priors we specify are treated as adding one replicate representing an  $N(50, 15)$  distribution and the updated distributions are then compared, with the posterior probability of 0.9936 corresponding to the probability that group A has a lower mean than group B. In this example we also include a region of practical equivalence (ROPE) test<sup>4</sup> and see that the probability of the 89% Highest Density Interval (HDI) of the true difference of means being in  $[-5, 5]$  is 0. Not only does the interpretation of the Bayesian posterior probability address our scientific question as stated but this framework also provides robust ways to test equivalency and hypotheses about the effect size between groups.

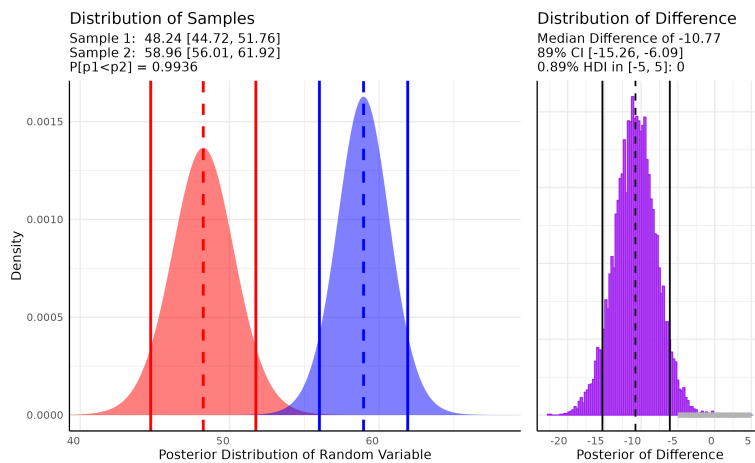


Figure 2: Bayesian Comparisons and Region of Practical Equivalence (ROPE) intervals provide much more information than standard frequentist tests for single or multi value traits.

through four `pcvr` functions. The first step in making growth models in this framework is to use `growthSS`, which specifies a self-starting non-linear growth model and returns a list of model components that are used by `fitGrowth` to fit the model. The fit model can be visualized using `growthPlot`, as shown in Figure 3. Finally, frequentist models can be tested against various nested versions of themselves using `testGrowth`, while more complex hypothesis tests are possible for Bayesian models using `brms::hypothesis`.<sup>10</sup> Using the `brms` backend, these models can also be combined into non-linear changepoint models, where changepoints are modeled per

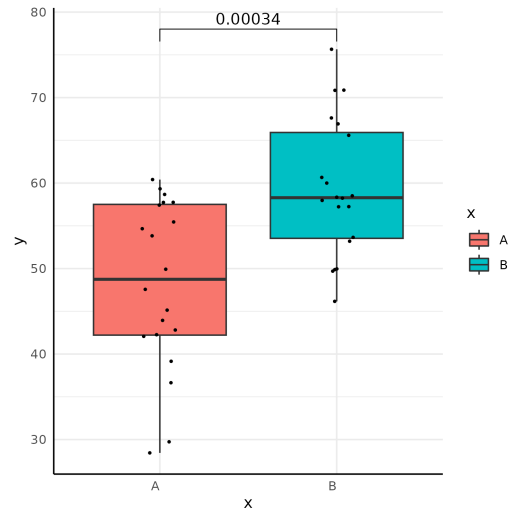


Figure 1: Single value, single timepoint traits can be compared using standard frequentist measures.

Longitudinal plant phenotype data presents several challenges due to auto-correlation within the data, often non-linear, and heteroskedastic (meaning its variance changes over time). Fitting non-linear models is typically more difficult than fitting linear models due to differences in how formulae are written and the need for starting values. In `pcvr` there are 8 parameterized growth curves<sup>5</sup> that can be fit using frequentist or Bayesian methods. For model fitting there are 4 supported back ends in R: `stats::nls`,<sup>2</sup> `quantreg::nlrq`,<sup>6</sup> `nlme::nlme`,<sup>7</sup> `brms::brm`.<sup>8</sup> Any of those 4 model fitting functions or `mgcv::gam`<sup>9</sup> can also fit generalized additive models. Specifying, visualizing, and testing complex growth models is greatly simplified

group. For instance, linear growth may persist for some time then give way to sigmoidal growth by passing "linear + logistic" to the `model` argument of `growthSS`. In this way very complex models of growth or stress recovery experiments can be specified easily in `pcvr`.

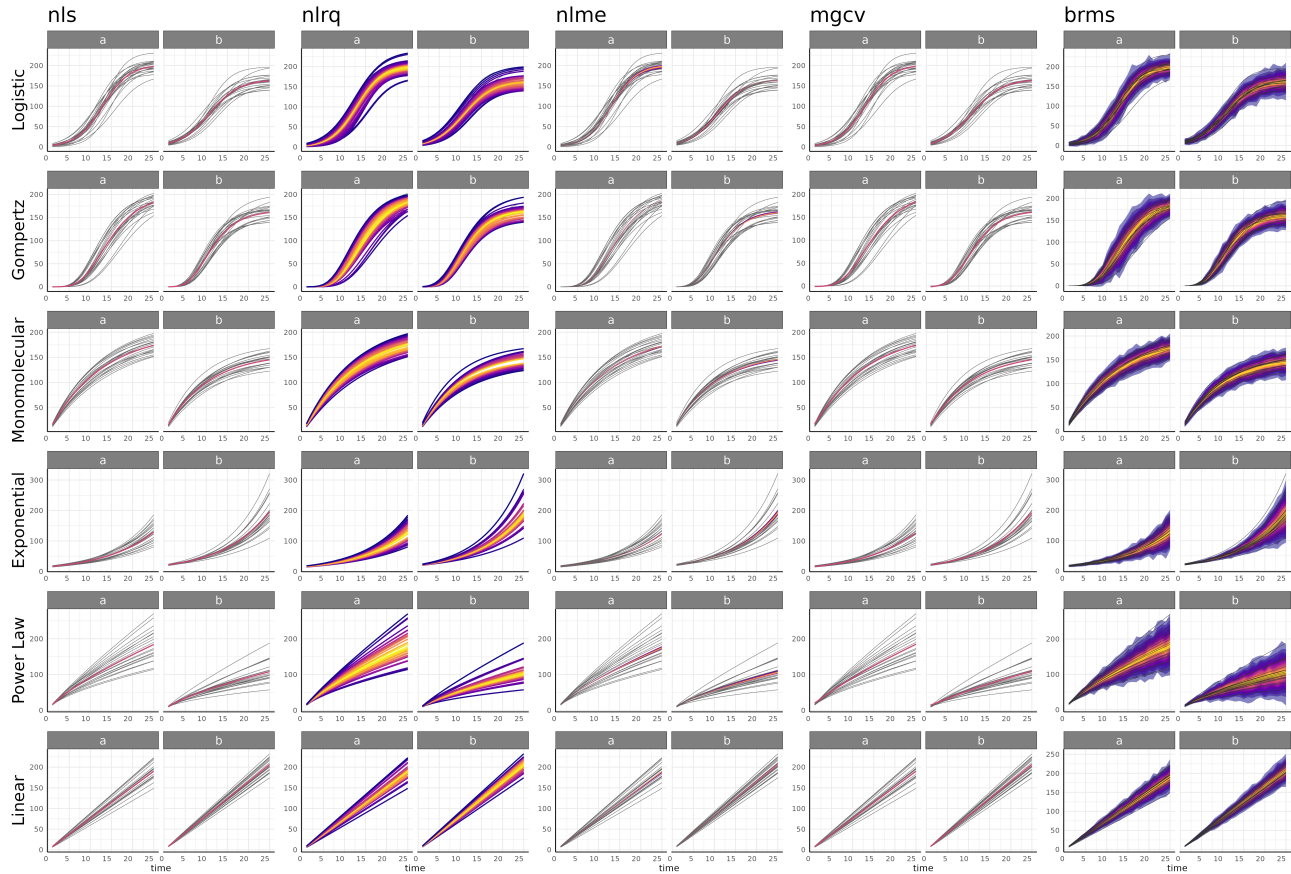


Figure 3: 6 parameterized self-starting growth models are shown as fit by 5 model fitting functions through `fitGrowth` and visualized by `growthPlot`.

Table 1: Each backend function has unique benefits, with `brms::brm` being the most versatile.

Backend Function	Non-Linearity	Autocorrelation	Heteroskedasticity	Changepoints
<code>stats::nls</code>	Yes	No	No	No
<code>mgcv::gam</code>	Yes	No	No	No
<code>quantreg::nlrq</code>	Yes	No	Yes *	No
<code>nlme::nlme</code>	Yes	Yes	Yes	No
<code>brms::brm</code>	Yes	Yes	Yes	Yes <sup>†</sup>

### 3. MULTI VALUE TRAIT ANALYSIS

Color and vegetative index data are returned from `PlantCV` as multi-value traits, meaning that the trait measured for each plant is described by a vector of numeric data. For example, the hue of pixels that comprise a plant may be represented as a histogram of each pixel's hue value. Multi-value traits present their own statistical

\*Changes in variance are expressed by different quantile fits.

<sup>†</sup>Double sigmoid models are available using each backend, but are not self-starting.

challenges since each observation contains a histogram. There are two main methods to analyze multi-value traits in **pcvr**: parametric Bayesian tests through **conjugate** and non-parametric analysis using Earth Mover’s Distance (EMD)<sup>11</sup> through **pcv.emd**. The **conjugate** function can take matrix input for samples, in which case the samples are assumed to be multi-value traits in wide format. Using **conjugate** region of practical equivalence (ROPE) tests and standard hypotheses can be conducted in the same manner as in single value trait analysis as shown in Figure 1.

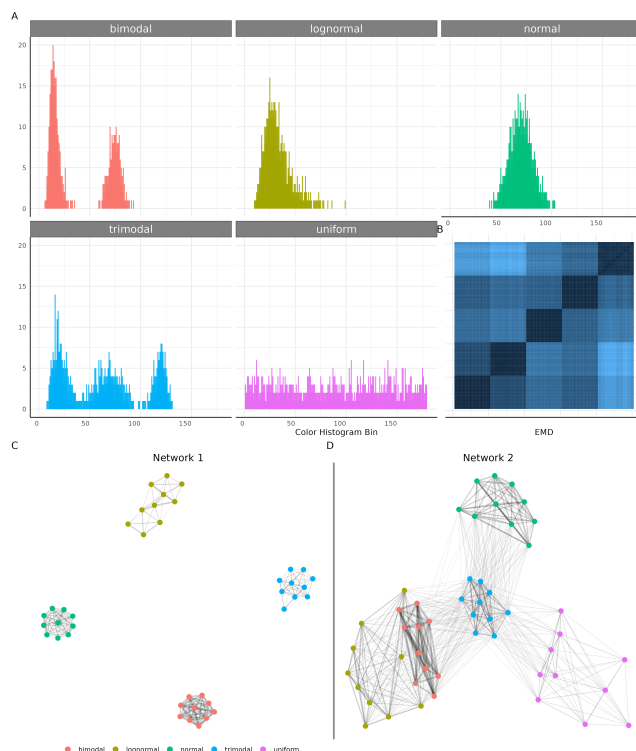


Figure 4: Non Parametric Multi Value Trait Analysis using EMD and network analysis easily shows distribution’s differences.

details and options for considering temporal trends in multi-value trait data are provided in both the [multi-value trait tutorial](#) and the Bellwether Phenotyping Facility vignette.

Color histograms can also be analyzed without considering them as probability densities by using EMD. EMD is a measure of work required to change one histogram into another. As a distance metric EMD is non-parametric, and is larger the more different two histograms are. The flexibility of EMD is demonstrated in Figure 4, where simulated image histograms from 5 distributions are compared with 10 replicates from each distribution. Data simulated from each distribution are shown in panel A and a heatmap of EMD values from pairwise comparisons between distributions as returned by **pcv.emd** in panel B. The heatmap shows that there are clear differences between the simulated distributions. Finally, EMD values are inverted to represent similarities and networks of those similarities are shown. Network 1 (panel C) is a network using only edges with a similarity score greater than 0.5 (scaled 0 to 1). This shows our distributions clustering together but has removed the uniform distribution since it does not contain edges stronger than our 0.5 cutoff. Network 2 (panel D) is a network using only edges stronger than the median edge strength and now we can see the relationships between these distributions, with bimodal and lognormal appearing closely related, trimodal being central, and uniform being more sparsely connected. Network analysis is used here for visual simplicity, but many distance-matrix methods are broadly reasonable here. Further

## 4. CONCLUSION

The **pcvr** package provides several useful tools to plant scientists, particularly those conducting high-throughput phenotyping experiments. The package is in active development and is presented with tutorials for several common sets of analyses. For a more complete example workflow please see the tutorials and [vignettes](#).

## 5. SOFTWARE AND DATA AVAILABILITY

All materials, data, and code used in this paper and in the tutorials/vignettes for **pcvr** are available on [GitHub](#).

## ACKNOWLEDGMENTS

Special thanks to Jeffrey Berry for help in learning and adopting more Bayesian statistics and for providing predecessors to the **conjugate** function. We also thank Leonardo Chavez and Joe Duenwald from the Bellwether Phenotyping Facility (RRID:SCR\_019049) at the Donald Danforth Plant Science Center for their expertise and support collecting Phenotyping Core Facility data used in **pcvr** vignettes. Finally we thank the Gehan lab for helping find bugs and providing interesting use cases.

## REFERENCES

- [1] Fahlgren, N., Feldman, M., Gehan, M., Wilson, M., Shyu, C., Bryant, D., Hill, S., McEntee, C., Warnasooriya, S., Kumar, I., Ficor, T., Turnipseed, S., Gilbert, K., Brutnell, T., Carrington, J., Mockler, T., and Baxter, I., “A versatile phenotyping system and analytics platform reveals diverse temporal responses to water availability in setaria,” *Molecular Plant* **8**(10), 1520–1535 (2015).
- [2] R Core Team, *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria (2013).
- [3] Fink, D., “A compendium of conjugate priors,” tech. rep., Montana State University (1997).
- [4] Kruschke, J., “Bayesian estimation supersedes the t test,” *Journal of Experimental Psychology: General* **142**(2), 573–603 (2013).
- [5] Paine, C. E. T., Marthens, T. R., Vogt, D. R., Purves, D., Rees, M., Hector, A., and Turnbull, L. A., “How to fit nonlinear plant growth models and calculate growth rates: An update for ecologists,” *Methods in Ecology and Evolution* **3**(2), 245–256 (2012).
- [6] Koenker, R., *quantreg: Quantile Regression* (2023). R package version 5.96.
- [7] Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D., and R Core Team, *nlme: Linear and Nonlinear Mixed Effects Models* (2022). R package version 3.1-155.
- [8] Bürkner, P.-C., “Bayesian item response modeling in R with brms and Stan,” *Journal of Statistical Software* **100**(5), 1–54 (2021).
- [9] Wood, S. N., “Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models,” *Journal of the Royal Statistical Society (B)* **73**(1), 3–36 (2011).
- [10] Bürkner, P.-C., “Advanced Bayesian multilevel modeling with the R package brms,” *The R Journal* **10**(1), 395–411 (2018).
- [11] Rubner, Y., Tomasi, C., and Guibas, L., “A metric for distributions with applications to image databases,” in *Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271)*, 59–66 (1998).