

ggplot2 compatible Quantile-quantile plots in R

by Alexandre Almeida, Adam Loy, Heike Hofmann

Abstract An abstract of less than 150 words.

TODO:

- Abstract
- Give examples
 - Heike: BRFSS example
- Conclusion

Background

Univariate distributional assessment is a common thread throughout statistical analyses during both the exploratory and confirmatory stages. When we begin exploring a new data set we often consider the distribution of individual variables before moving on to explore multivariate relationships. After a model has been fit to a data set, we must assess whether the distributional assumptions made are reasonable, and if they are not, then we must understand the impact this has on the conclusions of the model. Graphics provide arguably the most common way to carry out these univariate assessments. While there are many plots that can be used for distribution exploration and assessment, a quantile-quantile (Q-Q) plot (Wilk and Gnanadesikan, 1968) is one of the most common plots used.

Q-Q plots compare two distributions by matching a common set of quantiles. To compare a sample, y_1, y_2, \dots, y_n to a theoretical distribution, a Q-Q plot is simply a scatterplot of the sample quantiles, $y_{(i)}$, against the corresponding quantiles from the theoretical distribution, $F^{-1}(F_n(y_i))$. If the empirical distribution is consistent with the theoretical distribution, then the Q-Q plot will be linear. For example, Figure 1 shows two Q-Q plots: the left plot compares a sample drawn from the lognormal distribution to the lognormal distribution, while the right plot compares a sample drawn from the lognormal distribution to the normal distribution. As expected, the lognormal Q-Q plot is approximately linear as the data and model are in agreement, while the normal Q-Q plot is curved, indicating disagreement between the data and the model.

Additional graphical elements are often added to Q-Q plots in order to aid in distributional assessment. A reference line is often added to a Q-Q plot to assist the detection of departures from normality. This line is often drawn by connecting the first and third quartiles. Pointwise or simultaneous confidence bands are also frequently built around the reference line to display the expected degree of sampling error for the proposed model, so that minor deviations from the reference

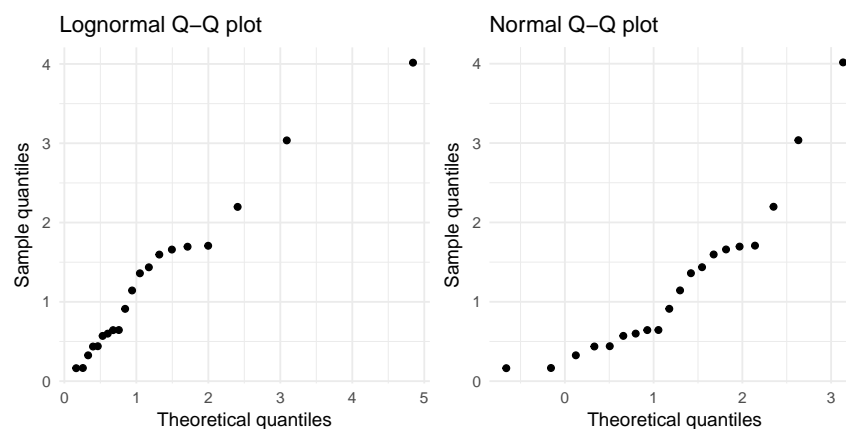


Figure 1: The left plot compares a sample drawn from the lognormal distribution to the lognormal distribution, while the right plot compares a sample drawn from the lognormal distribution to the normal distribution. The curvature in the normal Q-Q plot highlights the disagreement between the data and the model.

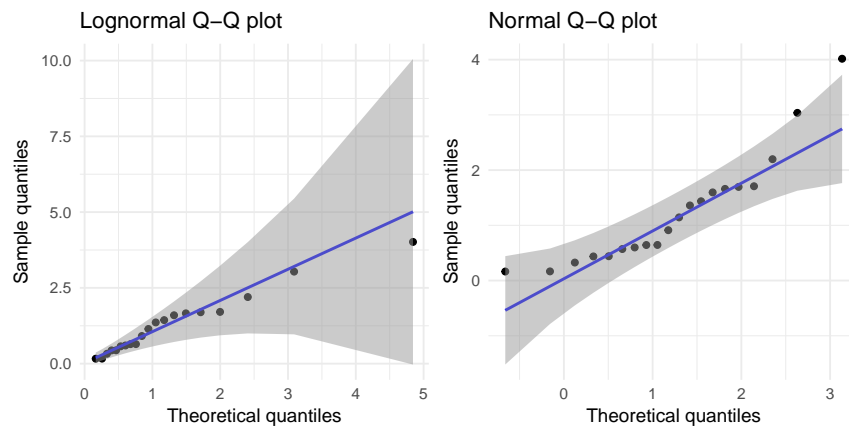


Figure 2: Adding reference lines and 95% pointwise confidence bands to the Q-Q plots in Figure 1.

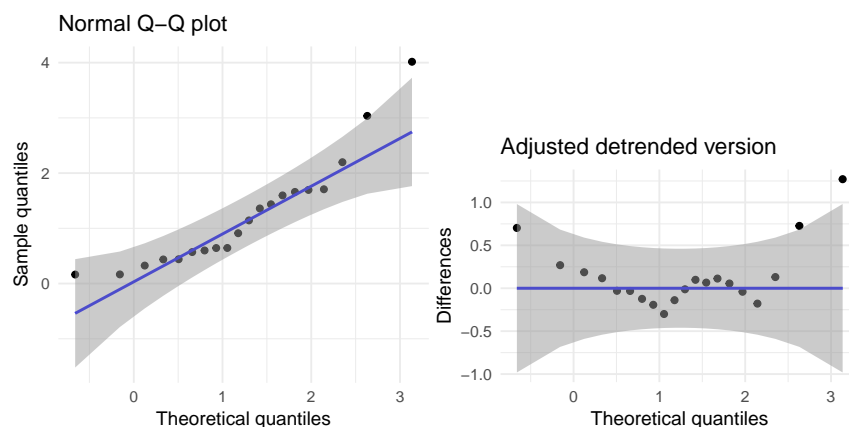


Figure 3: The left plot displays a traditional normal Q-Q plot for data simulated from a lognormal distribution. The right plot displays an adjusted detrended Q-Q plot of the same data, created by plotting the differences between the sample quantiles and the proposed model on the y -axis.

line are not over-interpreted. Figure 2 adds such reference lines and 95% pointwise confidence bands to the Q-Q plots in Figure 1.

Different orientations of Q-Q plots have also been proposed, most notably the de-trended Q-Q plot. To detrend a Q-Q plot, the y -axis is changed to show the difference between the observed quantile and the reference line. Consequently, the line representing the agreement with the theoretical distribution is the x -axis. Loy et al. (2016) find that detrended Q-Q plots are more powerful than other designs, so long as x - and y -axes to show the same range. Adjusting the aspect ratio in this way ensures that distances in the x - and y -directions are on the same scale. This Q-Q plot design is called an *adjusted-detrended Q-Q plot*. If the range of the axes are not adjusted, then this results an ordinary detrended Q-Q plot, which was found to have lower power than the standard Q-Q plot in some situations (Loy et al., 2016). Figure 3 displays the normal Q-Q plot from Figure 2 along with its adjusted detrended version.

Various implementations of Q-Q plots exist in R. Normal Q-Q plots, where a sample is compared to the standard normal distribution, are implemented using `qqplot` and `qqline` in **base** graphics (R Core Team, 2012). **lattice** provides a general framework for Q-Q plots in the `qqmath` function, allowing one to compare a sample to any theoretical distribution by specifying the appropriate quantile function (Sarkar, 2008). `qqPlot` in the **car** package also allows for the assessment of non-normal distributions and adds pointwise confidence bands via normal theory or the parametric bootstrap (Fox and Weisberg, 2011). **ggplot2** provides `geom_qq` and `geom_qq_line`, enabling the creation of Q-Q plots with a reference line, much like those created using `qqmath` (Wickham, 2016). None of these general use packages allow for easy construction of de-trended Q-Q plots.

qqplotr extends **ggplot2** to provide a complete implementation of Q-Q plots. The package allows for quick construction of all Q-Q plot designs without sacrificing the flexibility of the **ggplot2** framework. In the remainder of this paper, we will introduce the plotting framework provided by **qqplotr** and provide multiple examples of how it can be used.

Implementing Q-Q plots in the **ggplot2** framework

With **qqplotr** we extend some of the original **ggplot2** Q-Q plot functionalities by permitting the drawing of Q-Q points, lines, and confidence bands. Our approach provides a **ggplot2** layering mechanism so that for each one of those plot elements we implemented a **ggplot2** “stat” (statistical transformation). In addition, we also implemented a **ggplot2** “geom” (geometrical object) specifically for the confidence bands. That geom permits a simpler way of handling graphical parameters, which will become clearer in the [Examples](#) section.

The Q-Q plot functions are divided into three statistical transformations. Below, we describe each one of those functions and also give an overview of its parameters.

stat_qq_point

This is a modified version of `stat_qq/geom_qq` from **ggplot2** that plots the sample quantiles versus the theoretical quantiles (as in [Figure 1](#)). The novelty of this implementation is an option to detrend the plotted points (see [Background](#)). Note that all other implemented functions in the **qqplotr** package also allow for this option. Below we present a complete call to `stat_qq_point` and highlight its default parameters values:

```
stat_qq_point(
  data = NULL,
  mapping = NULL,
  geom = "point",
  position = "identity",
  na.rm = TRUE,
  show.legend = NA,
  inherit.aes = TRUE,
  distribution = "norm",
  dparams = list(),
  detrend = FALSE,
  qtype = 7,
  qprobs = c(0.25, 0.75),
  ...
)
```

- Parameters such as `data`, `mapping`, `geom`, `position`, `na.rm`, `show.legend`, and `inherit.aes` are specific for **ggplot2** implementations.
- `distribution` is a character string that sets the theoretical probability distribution. Here, we followed the nomenclature from the **stats** package. You must not provide the full distribution function name (e.g., “dnorm”). Instead, just provide its suffix (e.g., “norm”). If you wish to provide a custom distribution (as we will do in [Section 2.0.4](#)), you may do so by first creating the density (PDF), distribution (CDF), quantile, and random functions (also by following **stats** package nomenclature, i.e., for “custom”, you must provide the `dcustom`, `pcustom`, `qcustom`, and `rcustom` functions).
- `dparams` is a named list to be provided alongside with the previously chosen distribution. By default, MLEs are used for the distributional parameters. By manually providing any of the distributional parameters, the MLE is automatically turned off. Please note that MLEs are currently only supported for distributions available in the **stats** package; if a custom distribution is used in *distribution* *all* of its parameters have to be provided in `dparams`.
- `detrend` is a boolean that controls whether the points should be detrended (as shown in [Figure 3](#)). More on that in [Section 2.0.4](#).
- `qtype` and `qprobs` are only used when `detrend = TRUE`. These parameters are passed on to the `type` and `probs` parameters of quantile function from the **stats** package. The quantile function is used to define the quantiles where the reference line (i.e., the proposed model) shall intercept.

stat_qq_line

This statistical transformation draws a reference line based on the sample data quantiles. Note that the `stat_qq_line` call below does not present any additional parameters when compared to `stat_qq_point`.

```
stat_qq_line(
  data = NULL,
```

```

mapping = NULL,
geom = "path",
position = "identity",
na.rm = TRUE,
show.legend = NA,
inherit.aes = TRUE,
distribution = "norm",
dparams = list(),
detrend = FALSE,
qtype = 7,
qprobs = c(0.25, 0.75),
...
)

```

By default this line is drawn through two points, corresponding to the 25 and 75 percentile of the distributions on x (theoretical distribution) and y (empirical distribution). While we are using MLE to identify the exact distribution to compare against, these estimates are only used in determining the x -coordinates of these points. For the y coordinates 1st and 3rd quartile are used by default for a more robust estimation. The difference between observed and theoretical distribution therefore presents itself in the difference between the Q-Q line and the identity line. Using robust estimates for the Q-Q line is of particular advantage for small samples (Loy et al., 2016).

stat_qq_band

Draws confidence bands around the reference line using one of three methods: a normal approximation, the parametric bootstrap, or the tail-sensitive procedure. Below, we present the call to `stat_qq_band` and describe its specific parameters:

```

stat_qq_band(
  data = NULL,
  mapping = NULL,
  geom = "qq_band",
  position = "identity",
  show.legend = NA,
  inherit.aes = TRUE,
  na.rm = TRUE,
  distribution = "norm",
  dparams = list(),
  detrend = FALSE,
  qtype = 7,
  qprobs = c(0.25, 0.75),
  bandType = "normal",
  B = 1000,
  conf = 0.95,
  mu = NULL,
  sigma = NULL,
  ...
)

```

- `bandType` is a character string that controls the method to be used when constructing the confidence bands:
 - **Normal:** Specifying `bandType = "normal"` constructs pointwise confidence bands based on the normal approximation to the distribution of the order statistics. For example, an approximate 95% confidence interval for the i th order statistic is $\hat{X}_{(i)} \pm \Phi^{-1}(.975) \cdot SE(X_{(i)})$, where $\hat{X}_{(i)}$ denotes the value along the fitted line, $\Phi^{-1}(\cdot)$ denotes the quantile function for the standard normal distribution, and $SE(X_{(i)})$ is the standard error of the i th order statistic.
 - **Bootstrap:** Specifying `bandType = "bs"` constructs pointwise confidence bands using percentile confidence intervals from the parametric bootstrap.
 - **Tail-sensitive:** Specifying `bandType = "ts"` constructs the simulation-based tail-sensitive simultaneous confidence bands proposed by Aldor-Noiman et al. (2013).
- `B` is an integer that represents the number of bootstrap replicates (if `bandType = "bs"`) or the number of simulated samples (if `bandType = "ts"`).

- `conf` is a numerical variable that controls the confidence level of the bands, which must be a value between 0 and 1.
- `mu` and `sigma` are only used when `bandType = "ts"`. They represent the center and scale distributional parameters, respectively, which are used to construct the simulated tail-sensitive confidence bands. If any of the parameters is provided with `NULL`, then both the center and scale parameters will be estimated using robust estimates via the `robustbase` package.

facetting and qqplotr

placeholder for now: discuss groups and facets

Examples

In this section, we demonstrate the capabilities of the `qqplotr` package. We start by loading the package:

```
# also loads ggplot2
library(qqplotr)
```

BRFSS example

The Center for Disease Control and Prevention runs an annual telephone survey, the Behavioral Risk Factor Surveillance System (BRFSS), to keep track of the US populations' 'health-related risk behaviors, chronic health conditions, and use of preventive services'.

Close to half a million interviews are conducted each year. Here, we are focussing on the 2012 responses for Iowa. 7166 responses were gathered across 359 questions and derived variables. Among these, are participants' heights and weights, which we are going to assess in more detail.

Figure 4 shows two Q-Q plots side by side. For each of the plots, a sample of 200 men and 200 women is drawn from the overall number of responses. On the left hand side, individuals' heights are plotted in a Q-Q plot comparing raw heights to a normal distribution. We see that the distributions for both men and women (colour) is showing horizontal steps: this indicates that the distributional assessment is heavily dominated by the discreteness in the data, as most survey participants responded to the question of their height only up to the closest inch. On the right hand side of Figure 4, we use jittering; this means that we add a random number generated from a random uniform distribution on ± 0.5 inch to the reported height, as shown in the code below:

```
p2 <- sample_ia %>% mutate(
  HTIN4.jitter = jitter(HTIN4, factor = 2)
) %>%
  ggplot(aes(sample = HTIN4.jitter, colour=Gender, fill=Gender)) +
  geom_abline(slope=1, intercept=0, colour = "grey40") +
  stat_qq_band(alpha = 0.3, dparams = list(mean=params$m, sd=params$s)) +
  stat_qq_line(dparams = list(mean=params$m, sd=params$s)) +
  stat_qq_point(dparams = list(mean=params$m, sd=params$s)) +
  customization +
  ylab("Jittered Height (in inch)")
```

By using jittering we diminish the effect that discreteness has on the distribution and brings the observed distribution much closer to a normal distribution. Note that separate normal distributions were fitted for each gender. Not surprisingly, the resulting distributions have different means (women are on average 6 inch shorter than men in this dataset). Interestingly, the slope of the two genders is similar, indicating that the same scale parameter fits both genders' distributions (the standard deviation of height in the data set is 2.97 inch for men and 2.91 inch for women, see Table 1). The dark line between the two groups is the identity line – representing the theoretical distribution each of these groups are compared to. This distribution is based on parameters estimated from the whole population (see Table 1 for numbers). While the mean is about half way between the gender means, we see from the higher slope of the line that in comparison to each group, the standard deviation of the height based on the whole population is larger.

Unlike respondents' heights, their weights do not seem to be normally distributed. Figure 5 shows again two Q-Q plots. The Q-Q plot on the left uses raw weights and compares to a normal distribution. From the deviation in the extremes we see that tails of the observed distribution are heavier than expected under a normal distribution. On the right, weights are log-transformed. We see that a normal

Table 1: Summary of Iowa’s residents’ heights and weights with corresponding standard deviations by gender and for the total population.

SEX	mean height (in inch)	sd (in inch)	mean log weight (in kg)	sd (in kg)
Male	70.55	2.97	9.10	0.20
Female	64.51	2.91	8.89	0.23
Total	66.99	4.18	8.98	0.24

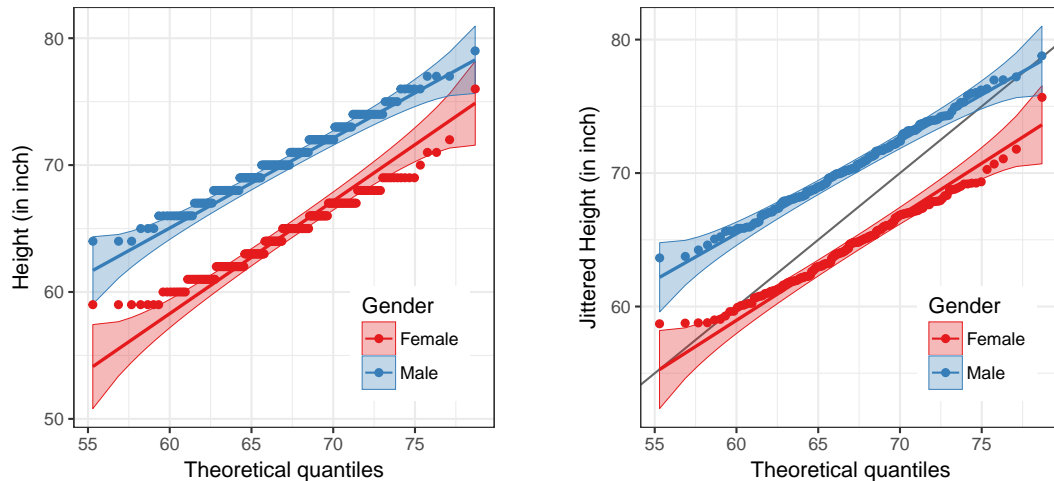


Figure 4: Sample (200 men and 200 women) of raw heights (left) and jittered heights (right). The distribution on the left is dominated by the discreteness of the data. On the right we see that except for some outliers an assumption of normality for people’s height is not completely absurd.

distribution for each of the genders shows –with the exceptions of a few extreme outliers– a reasonable fit.

Instead of transforming the observed values, we can change the theoretical distribution against which we compare. Figure 6 shows two Q-Q plots, one for each of the genders. A log-normal distribution is chosen as the theoretical distribution. By default, for each of the groups ML estimates are used:

```
p6 <- sample_ia %>%
  ggplot(aes(sample = WTKG3 / 100, colour = Gender, fill = Gender)) +
  geom_abline(colour="grey60") +
  stat_qq_band(distribution = "lnorm", alpha = 0.3) +
  stat_qq_line(distribution = "lnorm") +
  stat_qq_point(distribution = "lnorm") +
  scale_fill_brewer(palette = "Set1") +
  scale_colour_brewer(palette = "Set1") +
  customization + facet_grid(.~Gender)
```

Using user-provided distributions

Question: Should we move these sections above the BRFSS section, since they are a bit more straightforward?

Using the capabilities of **qqplotr** with the distributions implemented in the **stats** package is relatively straightforward, since the implementation allows you to specify the suffix (i.e. distribution and or abbreviation) via the `distribution` argument and the parameter estimates via `dparams` argument. However, there are times when the distributions in **stats** are not sufficient for the demands of the analysis. For example, there is no left-skewed distribution listed. User-coded distributions or distributions from other packages can be used with **qqplotr** as long as the distributions are defined following the conventions laid out in the **stats** package. Specifically, for some distribution there must be density/mass (`d` prefix), CDF (`p` prefix), quantile (`q` prefix), and simulation (`r` prefix) functions. In this section we illustrate the use of the smallest extreme value distribution (SEV).

To qualify for the Olympics in the men’s long jump in 2012, athletes had to either meet/exceed the

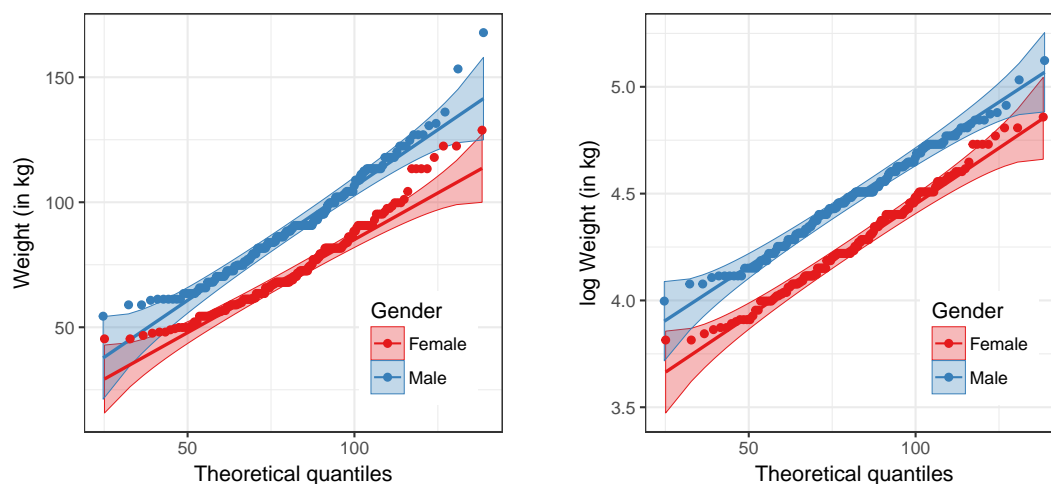


Figure 5: Sample (200 men and 200 women) of weights. Unlike people's height, weight seems to be heavily right skewed with some additional outliers on the extreme left (left plot). On the right, weight was log-transformed before its distribution is compared to a theoretical normal.

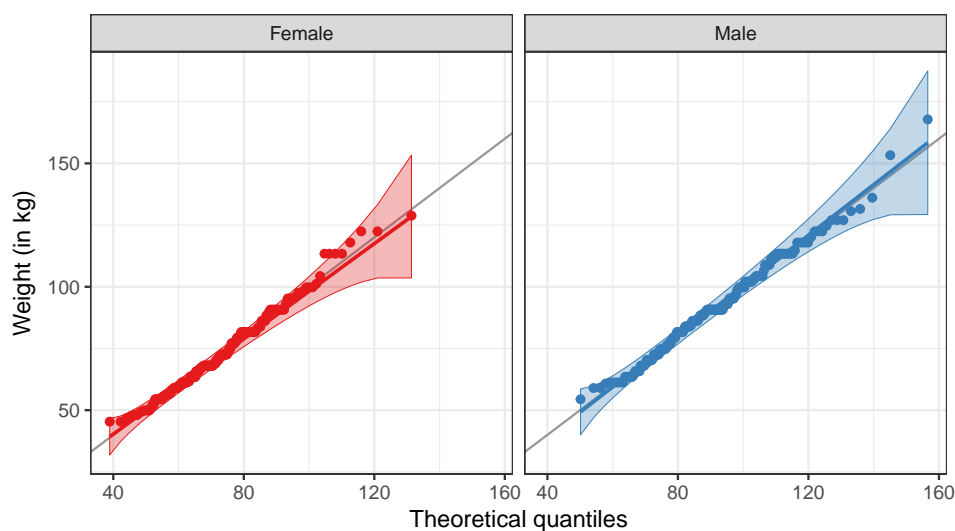


Figure 6: Sample (200 men and 200 women) of weights, the theoretical distribution is changed to a log normal. Shift and scale parameters are estimated separately for each of the genders before comparing distributions to a log-normal.

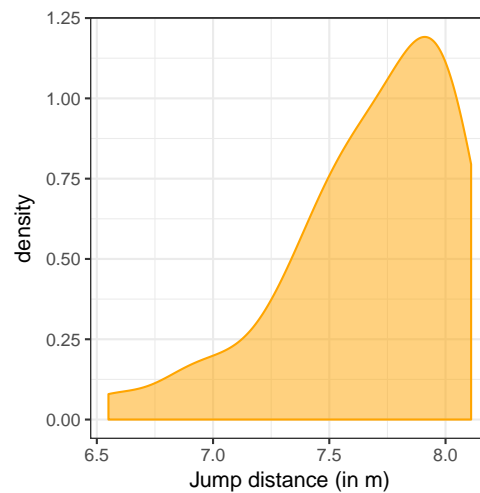


Figure 7: Density plot of the 2012 men's long jump qualifying round. The distances are clearly left skewed.

8.1 meter standard or place in the top twelve. During the qualification events, each athlete was able to jump three times, and their best (i.e. longest) jump is treated as the result. Figure 7 shows a density plot of the results, which are clearly left skewed.

In order to model the jump distances we must first define a left-skewed distribution. Below, we define the suite of distributional functions necessary to utilize the SEV distribution.

```
# CDF
psev <- function(q, mu = 0, sigma = 1) {
  z <- (q - mu) / sigma
  1 - exp(-exp(z))
}

# PDF
dsev <- function(x, mu = 0, sigma = 1) {
  z <- (x - mu) / sigma
  (1 / sigma) * exp(z - exp(z))
}

# Quantile function
qsev <- function(p, mu = 0, sigma = 1) {
  mu + log(-log(1 - p)) * sigma
}

# Simulation function
rsev <- function(n, mu = 0, sigma = 1) {
  qsev(runif(n), mu, sigma)
}
```

With the `*sev` distribution functions in hand, we can create a Q-Q plot to assess the appropriateness of the SEV model (Figure 8). The Q-Q plot show that the distances do not substantially deviate from the SEV model, so we have found an adequate representation of the distances.

```
ggplot(longjump, aes(sample = distance)) +
  stat_qq_band(distribution = "sev", dparams = list(mu = 0, sigma = 1), alpha = 0.3) +
  stat_qq_line(distribution = "sev", dparams = list(mu = 0, sigma = 1)) +
  stat_qq_point(distribution = "sev", dparams = list(mu = 0, sigma = 1)) +
  xlab("Theoretical quantiles") +
  ylab("Jump distance (in m)") +
  theme_bw()
```

Detrending Q-Q plots

AL: It seems natural to have a subsection here, even if we continue the example, to allow people

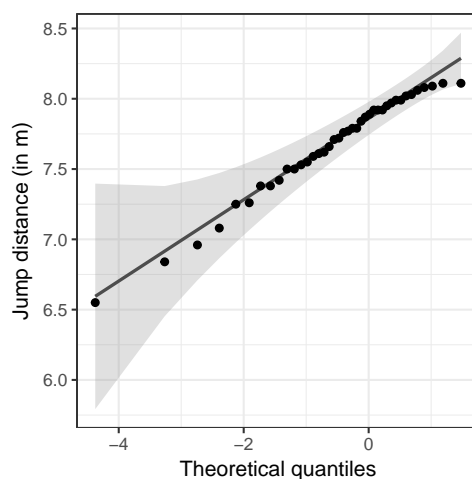


Figure 8: Q-Q plot comparing the long jump distances to the standard SEV distribution. The SEV distribution appears to adequately model the distances.

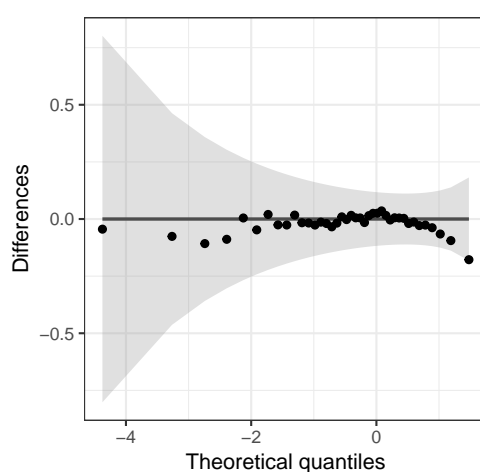


Figure 9: An adjusted detrended Q-Q plot assessing the appropriateness of the SEV distribution for the long jump data.

to skim the paper and find the topic they are most interested in. I tried to eliminate the technical discussion of the plot since it is in the background section, instead describing the code. What do you think?

To illustrate how to construct an adjusted detrended Q-Q plot using `qqplotr`, consider detrending the Q-Q plot in Figure 8. To obtain a detrended version of Figure 8, we must add the argument `detrend = TRUE` to `stat_qq_point`, `stat_qq_line`, and `stat_qq_band`. To adjust the aspect ratio to ensure that vertical and horizontal distances are on the same scale we further add `theme(aspect.ratio = 1)`.

```
ggplot(longjump, aes(sample = distance)) +
  stat_qq_band(distribution = "sev", alpha = 0.3, detrend = TRUE, dparams = c(mu = 0, sigma = 1)) +
  stat_qq_line(distribution = "sev", detrend = TRUE, dparams = c(mu = 0, sigma = 1)) +
  stat_qq_point(distribution = "sev", detrend = TRUE, dparams = c(mu = 0, sigma = 1)) +
  xlab("Theoretical quantiles") +
  ylab("Differences") +
  theme_bw() +
  theme(aspect.ratio = 1)
```

Q-Q plots with robust estimators

Add text

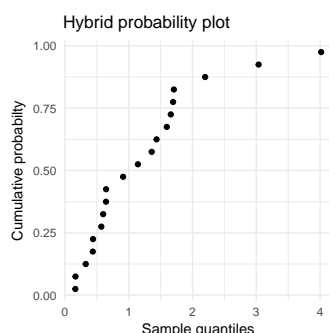


Figure 10: Illustrating different designs of probability plots.

Discussion

This paper presents the **qqplotr** package, an extension of **ggplot2** that implements Q-Q plots in both the standard and de-trended orientations, reference lines for Q-Q plots, and confidence bands for Q-Q plots. The examples illustrate how to create Q-Q plots for non-standard distributions found outside of the **stats**, de-trend Q-Q plots, and create Q-Q plots when data are grouped. Further, in the BRFSS example, we illustrated how jittering can be used in Q-Q plots to better compare discretized data to a continuous distribution.

Q-Q plots are member of the larger probability plotting family, and we are working to incorporate additional plots into **qqplotr**. . . NEED TO DISCUSS P-P PLOTS SOMEHOW

XXX P-P plots here Write this section once the rest of the paper is done.

Bibliography

- S. Aldor-Noiman, L. D. Brown, A. Buja, W. Rolke, and R. A. Stine. The power to see: A new graphical test of normality. *The American Statistician*, 67(4):249–260, 1 Nov. 2013. [p4]
- J. Fox and S. Weisberg. *An R Companion to Applied Regression*. Sage, Thousand Oaks CA, second edition, 2011. URL <http://socserv.socsci.mcmaster.ca/jfox/Books/Companion>. [p2]
- A. Loy, L. Follett, and H. Hofmann. Variations of Q–Q plots: The power of our eyes! *The American Statistician*, 70(2):202–214, 2 Apr. 2016. [p2, 4]
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2012. URL <http://www.R-project.org/>. ISBN 3-900051-07-0. [p2]
- D. Sarkar. *Lattice: Multivariate Data Visualization with R*. Springer, New York, 2008. URL <http://lmdvr.r-forge.r-project.org>. ISBN 978-0-387-75968-5. [p2]
- H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <http://ggplot2.org>. [p2]
- M. B. Wilk and R. Gnanadesikan. Probability plotting methods for the analysis of data. *Biometrika*, 55(1):1–17, Mar. 1968. [p1]

Acknowledgements

This work was partially funded by Google Summer of Code 2017.

Alexandre Almeida
University of Campinas
Institute of Computing
Campinas, Brazil 13083-852
almeida.xan@gmail.com

Adam Loy
Carleton College
Department of Mathematics and Statistics

Northfield, MN 55057
aloy@carleton.edu

Heike Hofmann
Iowa State University
Department of Statistics
Ames, IA 50011-1210
hofmann@iastate.edu