

Diagnostic plots — Distributional diagnostic plots

| | |
|--|---|
| Description | Quick start |
| Menu | Syntax |
| Options for <code>symplot</code>, <code>quantile</code>, and <code>qqplot</code> | Options for <code>qnorm</code> and <code>pnorm</code> |
| Options for <code>qchi</code> and <code>pchi</code> | Remarks and examples |
| Methods and formulas | Acknowledgments |
| References | Also see |

Description

`symplot` graphs a symmetry plot of *varname*.

`quantile` plots the ordered values of *varname* against the quantiles of a uniform distribution.

`qqplot` plots the quantiles of *varname*₁ against the quantiles of *varname*₂ (Q–Q plot).

`qnorm` plots the quantiles of *varname* against the quantiles of the normal distribution (Q–Q plot).

`pnorm` graphs a standardized normal probability plot (P–P plot).

`qchi` plots the quantiles of *varname* against the quantiles of a χ^2 distribution (Q–Q plot).

`pchi` graphs a χ^2 probability plot (P–P plot).

See [\[R\] regress postestimation diagnostic plots](#) for regression diagnostic plots and [\[R\] logistic postestimation](#) for logistic regression diagnostic plots.

Quick start

Symmetry plot for *v1*

```
symplot v1
```

Change marker color and size

```
symplot v1, mcolor(red) msize(large)
```

Plot ordered values of *v1* against quantiles of the uniform distribution

```
quantile v1
```

As above, but only for observations with *v2* greater than 5

```
quantile v1 if v2 > 5
```

Plot quantiles of *v1* against quantiles of *v2*

```
qqplot v1 v2
```

Change thickness of the reference line

```
qqplot v1 v2, rlopts(lwidth(thick))
```

Plot quantiles of *v1* against quantiles of the normal distribution

```
qnorm v1
```

Add grid lines

```
qnorm v1, grid
```

Standardized normal probability plot for v_1

```
pnorm v1
```

Change labels on the x and y axes

```
pnorm v1, xlabel(0(0.1)1) ylabel(0(0.1)1)
```

Plot quantiles of v_1 against quantiles of the χ_1^2 distribution

```
qchi v1
```

As above, but comparing with quantiles of the χ_2^2 distribution

```
qchi v1, df(2)
```

χ^2 probability plot for v_1

```
pchi v1
```

Add “ $\chi^2(1)$ P-P plot” to graph

```
pchi v1, title("{\&chi}{sup:2}(1) P-P plot")
```

Menu

symplot

Statistics > Summaries, tables, and tests > Distributional plots and tests > Symmetry plot

quantile

Statistics > Summaries, tables, and tests > Distributional plots and tests > Quantiles plot

qqplot

Statistics > Summaries, tables, and tests > Distributional plots and tests > Quantile–quantile plot

qnorm

Statistics > Summaries, tables, and tests > Distributional plots and tests > Normal quantile plot

pnorm

Statistics > Summaries, tables, and tests > Distributional plots and tests > Normal probability plot, standardized

qchi

Statistics > Summaries, tables, and tests > Distributional plots and tests > Chi-squared quantile plot

pchi

Statistics > Summaries, tables, and tests > Distributional plots and tests > Chi-squared probability plot

Syntax

Symmetry plot

```
symplot varname [if] [in] [ , options1 ]
```

Ordered values of *varname* against quantiles of uniform distribution

```
quantile varname [if] [in] [ , options1 ]
```

Quantiles of *varname*₁ against quantiles of *varname*₂

```
qqplot varname1 varname2 [if] [in] [ , options1 ]
```

Quantiles of *varname* against quantiles of normal distribution

```
qnorm varname [if] [in] [ , options2 ]
```

Standardized normal probability plot

```
pnorm varname [if] [in] [ , options2 ]
```

Quantiles of *varname* against quantiles of χ^2 distribution

```
qchi varname [if] [in] [ , options3 ]
```

χ^2 probability plot

```
pchi varname [if] [in] [ , options3 ]
```

| <i>options</i> ₁ | Description |
|---|---|
| Plot | |
| <i>marker_options</i> | change look of markers (color, size, etc.) |
| <i>marker_label_options</i> | add marker labels; change look or position |
| Reference line | |
| <u>rlopts</u> (<i>cline_options</i>) | affect rendition of the reference line |
| Add plots | |
| <i>addplot</i> (<i>plot</i>) | add other plots to the generated graph |
| Y axis, X axis, Titles, Legend, Overall | |
| <i>twoway_options</i> | any options other than by() documented in [G-3] <i>twoway_options</i> |

4 Diagnostic plots — Distributional diagnostic plots

| <i>options</i> ₂ | Description |
|---|--|
| Main | |
| <code>grid</code> | add grid lines |
| Plot | |
| <code>marker_options</code> | change look of markers (color, size, etc.) |
| <code>marker_label_options</code> | add marker labels; change look or position |
| Reference line | |
| <code>rlopts(cline_options)</code> | affect rendition of the reference line |
| Add plots | |
| <code>addplot(plot)</code> | add other plots to the generated graph |
| Y axis, X axis, Titles, Legend, Overall | |
| <code>twoway_options</code> | any options other than <code>by()</code> documented in [G-3] <code>twoway_options</code> |

| <i>options</i> ₃ | Description |
|---|--|
| Main | |
| <code>grid</code> | add grid lines |
| <code>df(#)</code> | degrees of freedom of χ^2 distribution; default is <code>df(1)</code> |
| Plot | |
| <code>marker_options</code> | change look of markers (color, size, etc.) |
| <code>marker_label_options</code> | add marker labels; change look or position |
| Reference line | |
| <code>rlopts(cline_options)</code> | affect rendition of the reference line |
| Add plots | |
| <code>addplot(plot)</code> | add other plots to the generated graph |
| Y axis, X axis, Titles, Legend, Overall | |
| <code>twoway_options</code> | any options other than <code>by()</code> documented in [G-3] <code>twoway_options</code> |

Options for `symplot`, `quantile`, and `qqplot`

Plot

`marker_options` affect the rendition of markers drawn at the plotted points, including their shape, size, color, and outline; see [G-3] `marker_options`.

`marker_label_options` specify if and how the markers are to be labeled; see [G-3] `marker_label_options`.

Reference line

`rlopts(cline_options)` affect the rendition of the reference line; see [G-3] `cline_options`.

Add plots

`addplot(plot)` provides a way to add other plots to the generated graph; see [G-3] [addplot_option](#).

Y axis, X axis, Titles, Legend, Overall

`twoway_options` are any of the options documented in [G-3] [twoway_options](#), excluding `by()`. These include options for titling the graph (see [G-3] [title_options](#)) and for saving the graph to disk (see [G-3] [saving_option](#)).

Options for `qnorm` and `pnorm`

Main

`grid` adds grid lines at the 0.05, 0.10, 0.25, 0.50, 0.75, 0.90, and 0.95 quantiles when specified with `qnorm`. With `pnorm`, `grid` is equivalent to `yline(.25,.5,.75) xline(.25,.5,.75)`.

Plot

`marker_options` affect the rendition of markers drawn at the plotted points, including their shape, size, color, and outline; see [G-3] [marker_options](#).

`marker_label_options` specify if and how the markers are to be labeled; see [G-3] [marker_label_options](#).

Reference line

`rlopts(cline_options)` affect the rendition of the reference line; see [G-3] [cline_options](#).

Add plots

`addplot(plot)` provides a way to add other plots to the generated graph; see [G-3] [addplot_option](#).

Y axis, X axis, Titles, Legend, Overall

`twoway_options` are any of the options documented in [G-3] [twoway_options](#), excluding `by()`. These include options for titling the graph (see [G-3] [title_options](#)) and for saving the graph to disk (see [G-3] [saving_option](#)).

Options for `qchi` and `pchi`

Main

`grid` adds grid lines at the 0.05, 0.10, 0.25, 0.50, 0.75, 0.90, and .95 quantiles when specified with `qchi`. With `pchi`, `grid` is equivalent to `yline(.25,.5,.75) xline(.25,.5,.75)`.

`df(#)` specifies the degrees of freedom of the χ^2 distribution. The default is `df(1)`.

Plot

`marker_options` affect the rendition of markers drawn at the plotted points, including their shape, size, color, and outline; see [G-3] [marker_options](#).

`marker_label_options` specify if and how the markers are to be labeled; see [G-3] [marker_label_options](#).

Reference line

`rlopts(cline_options)` affect the rendition of the reference line; see [G-3] [cline_options](#).

Add plots

`addplot(plot)` provides a way to add other plots to the generated graph; see [G-3] [addplot_option](#).

Y axis, X axis, Titles, Legend, Overall

twayway_options are any of the options documented in [G-3] [twayway_options](#), excluding `by()`. These include options for titling the graph (see [G-3] [title_options](#)) and for saving the graph to disk (see [G-3] [saving_option](#)).

Remarks and examples

[stata.com](http://www.stata.com)

Remarks are presented under the following headings:

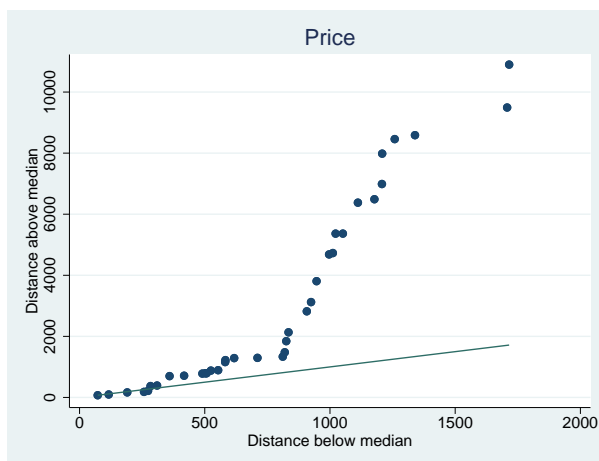
[symplot](#)
[quantile](#)
[qqplot](#)
[qnorm](#)
[pnorm](#)
[qchi](#)
[pchi](#)

symplot

► Example 1

We have data on 74 automobiles. To make a symmetry plot of the variable price, we type

```
. use https://www.stata-press.com/data/r17/auto
(1978 automobile data)
. symplot price
```



All points would lie along the reference line (defined as $y = x$) if car prices were symmetrically distributed. The points in this plot lie above the reference line, indicating that the distribution of car prices is skewed to the right—the most expensive cars are far more expensive than the least expensive cars are inexpensive.

The logic works as follows: a variable, z , is distributed symmetrically if

$$\text{median} - z_{(i)} = z_{(N+1-i)} - \text{median}$$

where $z_{(i)}$ indicates the i th-order statistic of z . `symplot` graphs $y_i = \text{median} - z_{(i)}$ versus $x_i = z_{(N+1-i)} - \text{median}$.

For instance, consider the largest and smallest values of `price` in the example above. The most expensive car costs \$15,906 and the least expensive, \$3,291. Let's compare these two cars with the typical car in the data and see how much more it costs to buy the most expensive car, and compare that with how much less it costs to buy the least expensive car. If the automobile price distribution is symmetric, the price differences would be the same.

Before we can make this comparison, we must agree on a definition for the word “typical”. Let's agree that “typical” means median. The price of the median car is \$5,006.50, so the most expensive car costs \$10,899.50 more than the median car, and the least expensive car costs \$1,715.50 less than the median car. We now have one piece of evidence that the car price distribution is not symmetric. We can repeat the experiment for the second-most-expensive car and the second-least-expensive car. We find that the second-most-expensive car costs \$9,494.50 more than the median car, and the second-least-expensive car costs \$1,707.50 less than the median car. We now have more evidence. We can continue doing this with the third most expensive and the third least expensive, and so on.

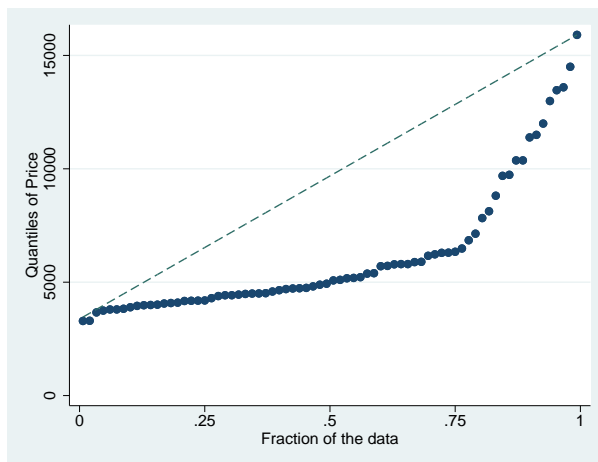
Once we have all of these numbers, we want to compare each pair and ask how similar, on average, they are. The easiest way to do that is to plot all the pairs.

quantile

▷ Example 2

We have data on the prices of 74 automobiles. To make a quantile plot of price, we type

```
. use https://www.stata-press.com/data/r17/auto, clear
(1978 automobile data)
. quantile price, rlopts(clpattern(dash))
```



We changed the pattern of the reference line by specifying `rlopts(clpattern(dash))`.

In a quantile plot, each value of the variable is plotted against the fraction of the data that have values less than that fraction. The diagonal line is a reference line. If automobile prices were rectangularly distributed, all the data would be plotted along the line. Because all the points are below the reference line, we know that the price distribution is skewed right.

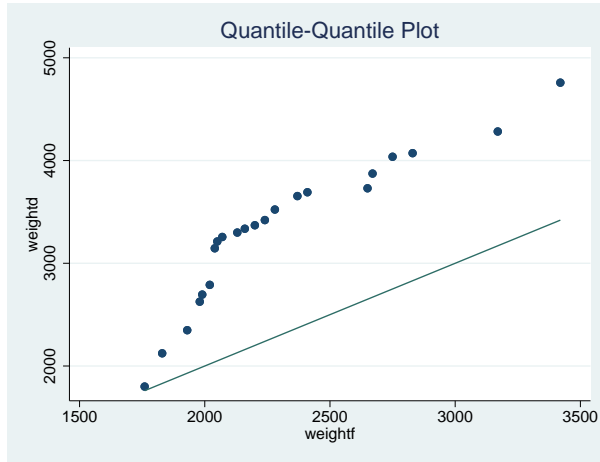
◀

qqplot

▷ Example 3

We have data on the weight and country of manufacture of 74 automobiles. We wish to compare the distributions of weights for domestic and foreign automobiles:

```
. use https://www.stata-press.com/data/r17/auto
(1978 automobile data)
. generate weightd=weight if !foreign
(22 missing values generated)
. generate weightf=weight if foreign
(52 missing values generated)
. qqplot weightd weightf
```

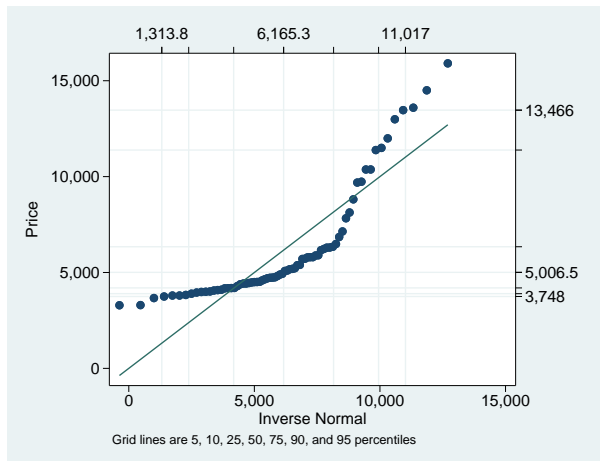



qnorm

▷ Example 4

Continuing with our price data on 74 automobiles, we now wish to compare the distribution of price with the normal distribution:

```
. qnorm price, grid ylabel(, angle(horizontal) axis(1))
> ylabel(, angle(horizontal) axis(2))
```



The result shows that the distributions are different.



□ Technical note

The idea behind `qnorm` is recommended strongly by Miller (1997): he calls it probit plotting. His recommendations from much practical experience should interest many users. “My recommendation for detecting nonnormality is *probit plotting*” (Miller 1997, 10). “If a deviation from normality cannot be spotted by eye on probit paper, it is not worth worrying about. I never use the Kolmogorov–Smirnov test (or one of its cousins) or the χ^2 test as a preliminary test of normality. They do not tell you how the sample is differing from normality, and I have a feeling they are more likely to detect irregularities in the middle of the distribution than in the tails” (Miller 1997, 13–14).

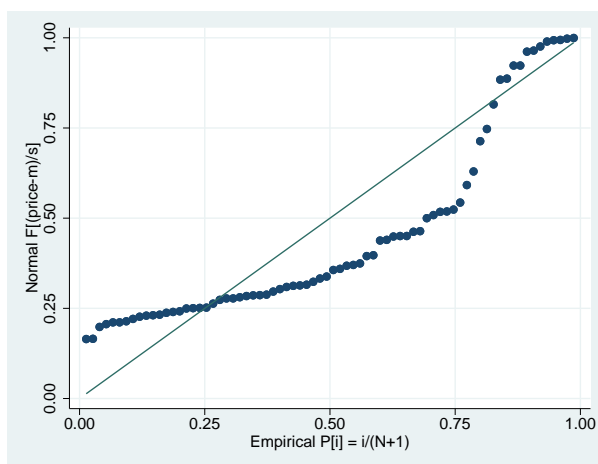
□

pnorm

▷ Example 5

Quantile–normal plots emphasize the tails of the distribution. Normal probability plots put the focus on the center of the distribution:

```
. pnorm price, grid
```



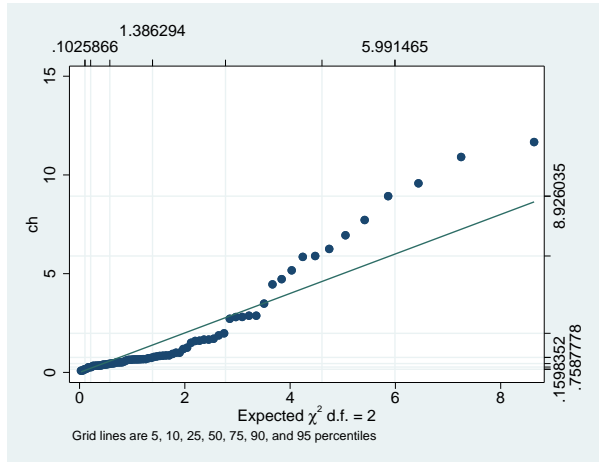
◀

qchi

▷ Example 6

Suppose that we want to examine the distribution of the sum of squares of `price` and `mpg`, standardized for their variances.

```
. egen c1 = std(price)
. egen c2 = std(mpg)
. generate ch = c1^2 + c2^2
. qchi ch, df(2) grid ylabel(, alt axis(2)) xlabel(, alt axis(2))
```



The quadratic form is clearly not χ^2 with 2 degrees of freedom.

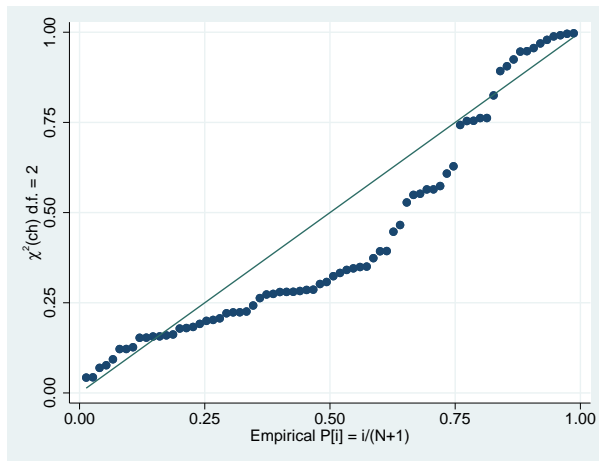


pchi

Example 7

We can focus on the center of the distribution by doing a probability plot:

```
. pchi ch, df(2) grid
```



Methods and formulas

Let $x_{(1)}, x_{(2)}, \dots, x_{(N)}$ be the data sorted in ascending order.

If a continuous variable, x , has a cumulative distribution function $F(x) = P(X \leq x) = p$, the quantiles x_{p_i} are such that $F(x_{p_i}) = p_i$. For example, if $p_i = 0.5$, then $x_{0.5}$ is the median. When we plot data, the probabilities, p_i , are often referred to as plotting positions. There are many different conventions for choice of plotting positions, given $x_{(1)} \leq \dots \leq x_{(N)}$. Most belong to the family $(i - a)/(N - 2a + 1)$. $a = 0.5$ (suggested by Hazen) and $a = 0$ (suggested by Weibull) are popular choices.

For a wider discussion of the calculation of plotting positions, see [Cox \(2002\)](#).

`symplot` plots median $- x_{(i)}$ versus $x_{(N+1-i)} - \text{median}$.

`quantile` plots $x_{(i)}$ versus $(i - 0.5)/N$ (the Hazen position).

`qnorm` plots $x_{(i)}$ against $q_i \times \hat{\sigma} + \hat{\mu}$, where $q_i = \Phi^{-1}(p_i)$, Φ is the cumulative normal distribution, $p_i = i/(N + 1)$ (the Weibull position), $\hat{\sigma}$ is the standard deviation, and $\hat{\mu}$ is the mean of the data.

`pnorm` plots $\Phi\{(x_i - \hat{\mu})/\hat{\sigma}\}$ versus $p_i = i/(N + 1)$, where $\hat{\mu}$ is the mean of the data and $\hat{\sigma}$ is the standard deviation.

`qchi` and `pchi` are similar to `qnorm` and `pnorm`; the cumulative χ^2 distribution is used in place of the cumulative normal distribution.

`qqplot` is just a two-way scatterplot of one variable against the other after both variables have been sorted into ascending order, and both variables have the same number of nonmissing observations. If the variables have unequal numbers of nonmissing observations, interpolated values of the variable with more data are plotted against the variable with fewer data.

Ramanathan Gnanadesikan (1932–2015) was born in Madras. He obtained degrees from the Universities of Madras and North Carolina. He worked in industry at Procter & Gamble, Bell Labs, and Bellcore, as well as in universities, retiring from Rutgers in 1998. Among many contributions to statistics, he is especially well known for work on probability plotting, robustness, outlier detection, clustering, classification, and pattern recognition.

Martin Bradbury Wilk (1922–2013) was born in Montreal. He obtained degrees in chemical engineering and statistics from McGill and Iowa State Universities. After holding several statistics-related posts in industry and at universities (including periods at Princeton, Bell Labs, and Rutgers), Wilk was appointed Chief Statistician at Statistics Canada (1980–1986). He is especially well known for his work with Gnanadesikan on probability plotting and with Shapiro on tests for normality.

Acknowledgments

We thank Peter A. Lachenbruch, Emeritus Appointment, Biostatistics, College of Public Health and Human Sciences, Oregon State University for writing the original versions of `qchi` and `pchi`. Patrick Royston of the MRC Clinical Trials Unit, London, and coauthor of the Stata Press book *Flexible Parametric Survival Analysis Using Stata: Beyond the Cox Model* also published a similar command in the *Stata Technical Bulletin* (Royston 1996).

References

- Chambers, J. M., W. S. Cleveland, B. Kleiner, and P. A. Tukey. 1983. *Graphical Methods for Data Analysis*. Belmont, CA: Wadsworth.
- Cox, N. J. 2002. Speaking Stata: On getting functions to do the work. *Stata Journal* 2: 411–427.
- . 2004a. Speaking Stata: Graphing distributions. *Stata Journal* 4: 66–88.
- . 2004b. [gr42_2](#): Software update: Quantile plots, generalized. *Stata Journal* 4: 97.
- . 2005a. Speaking Stata: Density probability plots. *Stata Journal* 5: 259–273.
- . 2005b. Speaking Stata: The protean quantile plot. *Stata Journal* 5: 442–460.
- . 2005c. Speaking Stata: Smoothing in various directions. *Stata Journal* 5: 574–593.
- . 2007. Stata tip 47: Quantile–quantile plots without programming. *Stata Journal* 7: 275–279.
- . 2012. Speaking Stata: Axis practice, or what goes where on a graph. *Stata Journal* 12: 549–561.
- Daniel, C., and F. S. Wood. 1980. *Fitting Equations to Data: Computer Analysis of Multifactor Data*. 2nd ed. New York: Wiley.
- Gan, F. F., K. J. Koehler, and J. C. Thompson. 1991. Probability plots and distribution curves for assessing the fit of probability models. *American Statistician* 45: 14–21. <https://doi.org/10.2307/2685233>.
- Genest, C., and G. J. Brackstone. 2013. Obituary: Martin B. Wilk, 1922–2013. *IMS Bulletin* 42(4): 7–8.
- Hoaglin, D. C. 1985. Using quantiles to study shape. In *Exploring Data Tables, Trends, and Shapes*, ed. D. C. Hoaglin, C. F. Mosteller, and J. W. Tukey, 417–460. New York: Wiley.
- Kettenring, J. R. 2001. A conversation with Ramanathan Gnanadesikan. *Statistical Science* 16: 295–309. <https://doi.org/10.1214/ss/1009213730>.
- Miller, R. G., Jr. 1997. *Beyond ANOVA: Basics of Applied Statistics*. London: Chapman & Hall.
- Nolan, D., and T. Speed. 2000. *Stat Labs: Mathematical Statistics Through Applications*. New York: Springer.
- Royston, P. 1996. [sg47](#): A plot and a test for the χ^2 distribution. *Stata Technical Bulletin* 29: 26–27. Reprinted in *Stata Technical Bulletin Reprints*, vol. 5, pp. 142–144. College Station, TX: Stata Press.
- Wilk, M. B., and R. Gnanadesikan. 1968. Probability plotting methods for the analysis of data. *Biometrika* 55: 1–17. <https://doi.org/10.2307/2334448>.

Also see

- [R] [cumul](#) — Cumulative distribution
- [R] [kdensity](#) — Univariate kernel density estimation
- [R] [logistic postestimation](#) — Postestimation tools for logistic
- [R] [iv](#) — Letter-value displays
- [R] [regress postestimation diagnostic plots](#) — Postestimation plots for regress