

Normality and other assumptions

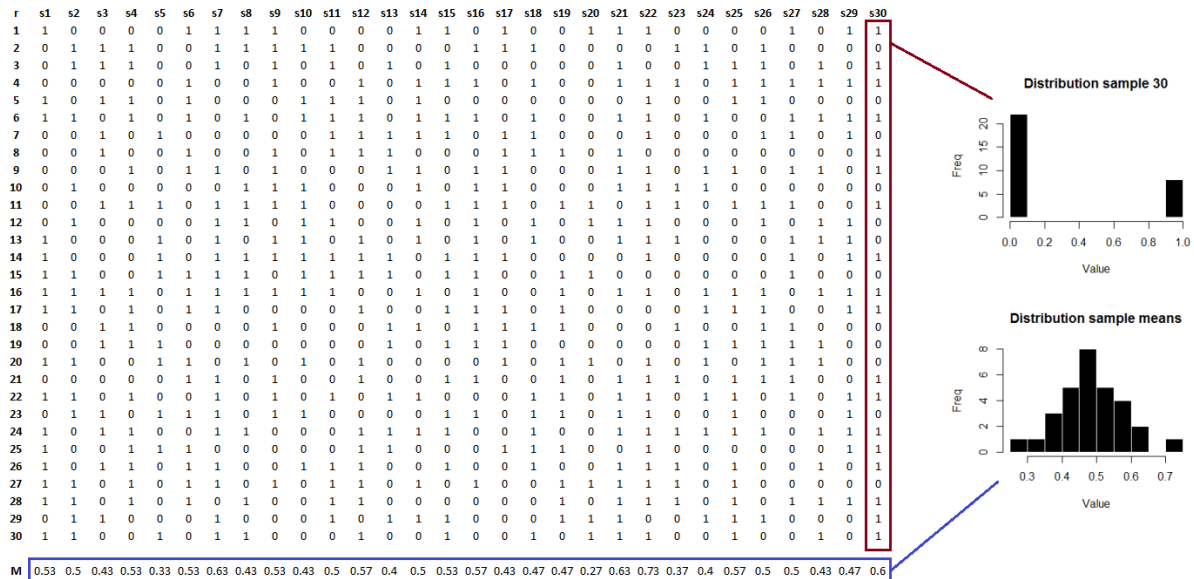
E-mail distributed on 11-12-2015

Dear colleagues,

I notice often that researchers are overly concerned with testing normality when checking statistical assumptions in an analysis. Today I would like to share two important reminders about normality.

1. Does it matter

In attachment I added an image that illustrates why normality is not necessarily important or influential in statistics. The picture shows 30 samples of data that have a non-normal distribution (Bernoulli in this case). However, as you can see, the averages of the samples are normally distributed!



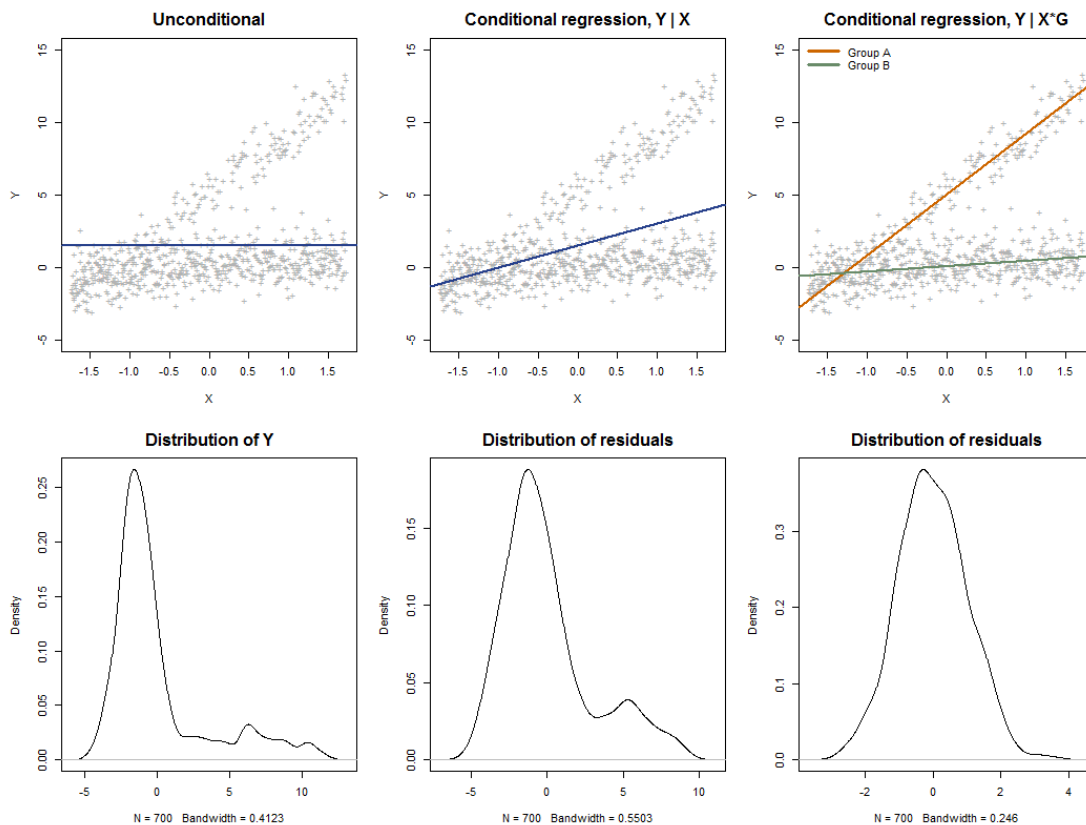
This is the result of a well-known theorem in statistics called the [central limit theorem](#), which says that the distribution of the sum (or average) of independent and identically distributed data will approximate a normal distribution, as the number of data increase to infinity. The image shows that the approximation is already quite good for N = 30!

It is important to remember that statistical tests are not conducted on raw data but always on some aggregated statistic (e.g., a mean, a variance). Distributional assumptions therefore relate to these statistics. If the raw data are already normally distributed then the statistic most likely will be too. However, as just shown, the statistic can be normally distributed without the raw data being so.

2. Conditional normality

Another point that is sometimes misunderstood is that normality is a *conditional* assumption. We assume that the dependent variable (DV) is normally distributed conditional on the effects in the model. For example, the classic independent-samples t-test assumes that the DV is normal within each group (conditional). It does not assume normality for the dependent across groups (unconditional).

In regression language, we assume that the dependent is normally distributed, given the independent variables in the model. This type of normality can be checked by inspecting the **residuals** of the model. In attachment I added graph that illustrates the difference between unconditional and conditional normality. In the graph, Y is conditionally normal (right panel) but unconditionally non-normal (left panel). I often notice that researchers only check for unconditional normality of Y. However, this would only be appropriate for the one-sample t-test!



The graph also illustrates that you should be careful with “correcting” non-normality through transformations (e.g., log-transform). Non-normal residuals may instead reflect a misspecification of

your underlying model. In the graph, Y is normal only when both X and G (and their interaction) are included as IVs in the model. You should explore such issues always before deciding to transform. For this reason it is also advised to screen your data visually (e.g., scatterplots, boxplots). Likewise, if you feel you must check for normality, it is recommended to make a visual inspection with a quantile-quantile plot (QQ-plot). This will give you more information than a statistical test (e.g., Kolmogorov-Smirnov).

Finally, I advise you to worry about statistical assumptions other than normality, such as **assumptions on variances and covariances** (e.g., homoscedasticity, sphericity), and to diagnose other types of problems of your fitted model, such as multicollinearity, linearity, and influential cases. Violations against these assumptions or issues tend to have a much more dramatic impact on statistical estimation and inference than violations of the normality assumption.

Have a nice weekend,
Ben

--

Ben Meuleman, Ph.D.
Statistical Assistant

Campus Biotech | CISA - University of Geneva
Chemin des Mines 9 | CH-1202 Genève
ben.meuleman@unige.ch | +41 (0)22 379 09 79