# Statistical tests for checking data normality

Since our knowledge of a population studied is generally incomplete, such as the population standard deviation $\sigma$, we generally assume that a particular distribution will adequately model the distribution for a particular variable. In chemical analysis, repeated analyses on a sample and random sampling from a population tend to fall under a normal distribution. However, like with any other statistical assumption, this should be checked and verified.

Unfortunately, although there are many statistical tests for departures from normality, few are sufficiently powerful to be useful on the relatively small data sets which are common in laboratory testing. If specific tests for normality are considered essential, perhaps the commonest ones are the Kolmogorov–Smirnov (K–S) test and the Shapiro–Wilk test.   These can be found in most statistical software. A more powerful modification of the K–S test is the Anderson–Darling test. Others include the chi-squared goodness–of–fit test and tests for significance skewness and kurtosis.   Many of these tests can be used to compare data against other distributions.

Another way of looking at this issue is to examine the presence of any outlier in the data set, because outliers represent a particular common departure from normal distributions in analytical data. We may also make use of the scatter plots and normal probability plots to review data sets and identify significant anomalies.

In this note, we would like to study a set of data which is very close to the critical zone to decide whether the assumption of normality still holds water.

For example, we obtained the following set of 12 analytical data:

| 42 | 46 | 39 | 45 | 42 | 48 |
|----|----|----|----|----|----|
| 47 | 43 | 31 | 42 | 45 | 48 |

The XLSTAT software gave the following statistical test results:

Data: Workbook = Book1 / Sheet = Sheet1 / Range = Sheet1!$A$1:$A$13 / 12 rows and 1 column

Significance level (%): 5

**Summary statistics (Data):**

| Variable | Observations | Obs. with missing data | Obs. without missing data | Min | Max | Mean | Std. deviation |
|----------|--------------|------------------------|---------------------------|-----|-----|------|----------------|
| Data | 12 | 0 | 12 | 31.000 | 48.000 | 43.167 | 4.726 |

**Shapiro-Wilk test (Data):**

| | |
|---|---|
| W | 0.849 |
| p-value (Two-tailed) | 0.036 |
| alpha | 0.05 |

**Test interpretation:**

H0: The variable from which the sample was extracted follows a Normal distribution.

Ha: The variable from which the sample was extracted does not follow a Normal distribution.

The risk to reject the null hypothesis H0 while it is true is lower than 3.61%.

**Anderson-Darling test (Data):**

| | |
|---|---|
| A2 | 0.621 |
| p-value (Two-tailed) | 0.080 |
| alpha | 0.05 |

**Test interpretation:**

H0: The variable from which the sample was extracted follows a Normal distribution.

Ha: The variable from which the sample was extracted does not follow a Normal distribution.

The risk to reject the null hypothesis H0 while it is true is 8.04%.

**Summary: of p-values**

**($\alpha$ = 0.05)**

| Variable\Test | Shapiro-Wilk | Anderson-Darling |
|---|---|---|
| Data | **0.036** | 0.080 |

The above example shows that whilst Shapiro-Wilk test rejected the null hypothesis Ho, indicating that the data set did not follow a normal distribution at $\alpha = 0.05$ level, the Anderson-Darling statistic test accepted the Ho, instead at the same level of confidence.   To solve this dilemma, we should take a closer look at the data set and check if there is any outlier figure that affects such normality testing.

Indeed, outlier tests such as Grubb's and Dixon's showed that the test result 31 was an outlier with 95% confidence and must be removed. Upon its removal from the data set, the balanced 11 test results indicated a trend of normal distribution by both statistic tests.

Hence, it is recommended to always carry out an outlier testing on a set of data before subjecting it for a normality checking.