

NAG Library Chapter Introduction
g01 – Simple Calculations on Statistical Data

Contents

1	Scope of the Chapter	2
2	Background to the Problems	2
2.1	Summary Statistics	2
2.2	Statistical Distribution Functions and Their Inverses	2
2.3	Testing for Normality and Other Distributions	3
2.4	Distribution of Quadratic Forms	3
2.5	Energy Loss Distributions	3
2.6	Vectorized Functions	4
3	Recommendations on Choice and Use of Available Functions	4
3.1	Working with Streamed or Extremely Large Datasets	7
4	Auxiliary Functions Associated with Library Function Arguments	7
5	Functions Withdrawn or Scheduled for Withdrawal	7
6	References	7

1 Scope of the Chapter

This chapter covers three topics:

- summary statistics
- statistical distribution functions and their inverses;
- testing for Normality and other distributions.

2 Background to the Problems

2.1 Summary Statistics

The summary statistics consist of two groups. The first group are those based on moments; for example mean, standard deviation, coefficient of skewness, and coefficient of kurtosis (sometimes called the ‘excess of kurtosis’, which has the value 0 for the Normal distribution). These statistics may be sensitive to extreme observations and some robust versions are available in Chapter g07. The second group of summary statistics are based on the order statistics, where the i th order statistic in a sample is the i th smallest observation in that sample. Examples of such statistics are minimum, maximum, median, hinges and quantiles.

2.2 Statistical Distribution Functions and Their Inverses

Statistical distributions are commonly used in three problems:

- evaluation of probabilities and expected frequencies for a distribution model;
- testing of hypotheses about the variables being observed;
- evaluation of confidence limits for parameters of fitted model, for example the mean of a Normal distribution.

Random variables can be either discrete (i.e., they can take only a limited number of values) or continuous (i.e., can take any value in a given range). However, for a large sample from a discrete distribution an approximation by a continuous distribution, usually the Normal distribution, can be used. Distributions commonly used as a model for discrete random variables are the binomial, hypergeometric, and Poisson distributions. The binomial distribution arises when there is a fixed probability of a selected outcome as in sampling with replacement, the hypergeometric distribution is used in sampling from a finite population without replacement, and the Poisson distribution is often used to model counts.

Distributions commonly used as a model for continuous random variables are the Normal, gamma, and beta distributions. The Normal is a symmetric distribution whereas the gamma is skewed and only appropriate for non-negative values. The beta is for variables in the range $[0, 1]$ and may take many different shapes. For circular data, the ‘equivalent’ to the Normal distribution is the von Mises distribution. The assumption of the Normal distribution leads to procedures for testing and interval estimation based on the χ^2 , F (variance ratio), and Student's t -distributions.

In the hypothesis testing situation, a statistic X with known distribution under the null hypothesis is evaluated, and the probability α of observing such a value or one more ‘extreme’ value is found. This probability (the significance) is usually then compared with a preassigned value (the significance level of the test), to decide whether the null hypothesis can be rejected in favour of an alternate hypothesis on the basis of the sample values. Many tests make use of those distributions derived from the Normal distribution as listed above, but for some tests specific distributions such as the Studentized range distribution and the distribution of the Durbin–Watson test have been derived. Nonparametric tests as given in Chapter g08, such as the Kolmogorov–Smirnov test, often use statistics with distributions specific to the test. The probability that the null hypothesis will be rejected when the simple alternate hypothesis is true (the power of the test) can be found from the noncentral distribution.

The confidence interval problem requires the inverse calculation. In other words, given a probability α , the value x is to be found, such that the probability that a value not exceeding x is observed is equal to α . A confidence interval of size $1 - 2\alpha$, for the quantity of interest, can then be computed as a function of x and the sample values.

The required statistics for either testing hypotheses or constructing confidence intervals can be computed with the aid of functions in this chapter, and Chapter g02 (for regression), Chapter g04 (for analysis of designed experiments), Chapter g13 (for time series), and Chapter e04 (for nonlinear least squares problems).

Pseudorandom numbers from many statistical distributions can be generated by functions in Chapter g05.

2.3 Testing for Normality and Other Distributions

Methods of checking that observations (or residuals from a model) come from a specified distribution, for example, the Normal distribution, are often based on order statistics. Graphical methods include the use of **probability plots**. These can be either $P - P$ plots (probability–probability plots), in which the empirical probabilities are plotted against the theoretical probabilities for the distribution, or $Q - Q$ plots (quantile–quantile plots), in which the sample points are plotted against the theoretical quantiles. $Q - Q$ plots are more common, partly because they are invariant to differences in scale and location. In either case if the observations come from the specified distribution then the plotted points should roughly lie on a straight line.

If y_i is the i th smallest observation from a sample of size n (i.e., the i th order statistic) then in a $Q - Q$ plot for a distribution with cumulative distribution function F , the value y_i is plotted against x_i , where $F(x_i) = (i - \alpha)/(n - 2\alpha + 1)$, a common value of α being $\frac{1}{2}$. For the Normal distribution, the $Q - Q$ plot is known as a Normal probability plot.

The values x_i used in $Q - Q$ plots can be regarded as approximations to the expected values of the order statistics. For a sample from a Normal distribution the expected values of the order statistics are known as **Normal scores** and for an exponential distribution they are known as **Savage scores**.

An alternative approach to probability plots are the more formal tests. A test for Normality is the Shapiro and Wilk's W Test, which uses Normal scores. Other tests are the χ^2 goodness-of-fit test and the Kolmogorov–Smirnov test; both can be found in Chapter g08.

2.4 Distribution of Quadratic Forms

Many test statistics for Normally distributed data lead to quadratic forms in Normal variables. If X is a n -dimensional Normal variable with mean μ and variance-covariance matrix Σ then for an n by n matrix A the quadratic form is

$$Q = X^T A X.$$

The distribution of Q depends on the relationship between A and Σ : if $A\Sigma$ is idempotent then the distribution of Q will be central or noncentral χ^2 depending on whether μ is zero.

The distribution of other statistics may be derived as the distribution of linear combinations of quadratic forms, for example the Durbin–Watson test statistic, or as ratios of quadratic forms. In some cases rather than the distribution of these functions of quadratic forms the values of the moments may be all that is required.

2.5 Energy Loss Distributions

An application of distributions in the field of high-energy physics where there is a requirement to model fluctuations in energy loss experienced by a particle passing through a layer of material. Three models are commonly used:

- (i) Gaussian (Normal) distribution;
- (ii) the Landau distribution;
- (iii) the Vavilov distribution.

Both the Landau and the Vavilov density functions can be defined in terms of a complex integral. The Vavilov distribution is the more general energy loss distribution with the Landau and Gaussian being suitable when the Vavilov parameter κ is less than 0.01 and greater than 10.0 respectively.

2.6 Vectorized Functions

A number of vectorized functions are included in this chapter. Unlike their scalar counterparts, which take a single set of parameters and perform a single function evaluation, these functions take vectors of parameters and perform multiple function evaluations in a single call. The input arrays to these vectorized functions are designed to allow maximum flexibility in the supply of the parameters by reusing, in a cyclic manner, elements of any arrays that are shorter than the number of functions to be evaluated, where the total number of functions evaluated is the size of the largest array.

To illustrate this we will consider `nag_prob_gamma_vector` (`g01sfc`), a vectorized version of `nag_gamma_dist` (`g01efc`), which calculates the probabilities for a gamma distribution. The gamma distribution has two parameters α and β therefore `nag_prob_gamma_vector` (`g01sfc`) has four input arrays, one indicating the tail required (**tail**), one giving the value of the gamma variate, g , whose probability is required (**g**), one for α (**a**) and one for β (**b**). The lengths of these arrays are **ltail**, **lg**, **la** and **lb** respectively.

For sake of argument, let's assume that **ltail** = 1, **lg** = 2, **la** = 3 and **lb** = 4, then $\max(\mathbf{ltail}, \mathbf{lg}, \mathbf{la}, \mathbf{lb}) = 4$ values will be returned. These four probabilities would be calculated using the following parameters:

i	Tail	g	α	β
1	tail [0]	g [0]	a [0]	b [0]
2	tail [0]	g [1]	a [1]	b [1]
3	tail [0]	g [0]	a [2]	b [2]
4	tail [0]	g [1]	a [0]	b [3]

3 Recommendations on Choice and Use of Available Functions

Descriptive statistics / Exploratory analysis,

summaries,

frequency / contingency table,

one variable `nag_frequency_table` (`g01aec`)

mean, variance, skewness, kurtosis (one variable),

combine summaries `nag_summary_stats_onevar_combine` (`g01auc`)

from frequency table `nag_summary_stats_freq` (`g01adc`)

from raw data `nag_summary_stats_onevar` (`g01atc`)

median, hinges / quartiles, minimum, maximum `nag_5pt_summary_stats` (`g01alc`)

quantiles,

approximate,

large data stream of fixed size `nag_approx_quantiles_fixed` (`g01anc`)

large data stream of unknown size `nag_approx_quantiles_arbitrary` (`g01apc`)

unordered vector `nag_double_quantiles` (`g01amc`)

rolling window,

mean, standard deviation (one variable) `nag_moving_average` (`g01wac`)

Distributions,

Beta,

central,

deviates,

scalar `nag_deviates_beta` (`g01fec`)

vectorized `nag_deviates_beta_vector` (`g01tec`)

probabilities and probability density function,

scalar `nag_prob_beta_dist` (`g01eec`)

vectorized `nag_prob_beta_vector` (`g01sec`)

non-central,

probabilities `nag_prob_non_central_beta_dist` (`g01gec`)

binomial,

distribution function,

scalar `nag_binomial_dist` (`g01bjc`)

vectorized `nag_prob_binomial_vector` (`g01sjc`)

Dickey–Fuller unit root test,

probabilities, `nag_prob_dickey_fuller_unit` (`g01ewc`)

Durbin–Watson statistic,	
probabilities	nag_prob_durbin_watson (g01epc)
energy loss distributions,	
Landau,	
density	nag_prob_density_landau (g01mtc)
derivative of density	nag_prob_der_landau (g01rtc)
distribution	nag_prob_landau (g01etc)
first moment	nag_moment_1_landau (g01ptc)
inverse distribution	nag_deviates_landau (g01ftc)
second moment	nag_moment_2_landau (g01qtc)
Vavilov,	
density	nag_prob_density_vavilov (g01muc)
distribution	nag_prob_vavilov (g01euc)
initialization	nag_init_vavilov (g01zuc)
<i>F</i> :	
central,	
deviates,	
scalar	nag_deviates_f_dist (g01fdc)
vectorized	nag_deviates_f_vector (g01tdc)
probabilities,	
scalar	nag_prob_f_dist (g01edc)
vectorized	nag_prob_f_vector (g01sdc)
non-central,	
probabilities	nag_prob_non_central_f_dist (g01gdc)
gamma,	
deviates,	
scalar	nag_deviates_gamma_dist (g01ffc)
vectorized	nag_deviates_gamma_vector (g01tfc)
probabilities,	
scalar	nag_gamma_dist (g01efc)
vectorized	nag_prob_gamma_vector (g01sfc)
probability density function,	
scalar	nag_gamma_pdf (g01kfc)
vectorized	nag_gamma_pdf_vector (g01kkc)
Hypergeometric,	
distribution function,	
scalar	nag_hypergeom_dist (g01blc)
vectorized	nag_prob_hypergeom_vector (g01slc)
Kolmogorov–Smirnov,	
probabilities,	
one-sample	nag_prob_1_sample_ks (g01eyc)
two-sample	nag_prob_2_sample_ks (g01ezc)
Normal,	
bivariate,	
probabilities	nag_bivariate_normal_dist (g01hac)
multivariate,	
probabilities	nag_multi_normal (g01hbc)
probability density function,	
vectorized	nag_multi_normal_pdf_vector (g01lbc)
quadratic forms,	
cumulants and moments	nag_moments_quad_form (g01nac)
moments of ratios	nag_moments_ratio_quad_forms (g01nbc)
univariate,	
deviates,	
scalar	nag_deviates_normal (g01fac)
vectorized	nag_deviates_normal_vector (g01tac)
probabilities,	
scalar	nag_prob_normal (g01eac)
vectorized	nag_prob_normal_vector (g01sac)

probability density function,	
scalar	nag_normal_pdf (g01kac)
vectorized	nag_normal_pdf_vector (g01kqc)
reciprocal of Mill's Ratio	nag_mills_ratio (g01mbc)
Shapiro and Wilk's test for Normality	nag_shapiro_wilk_test (g01ddc)
Poisson,	
distribution function,	
scalar	nag_poisson_dist (g01bkc)
vectorized	nag_prob_poisson_vector (g01skc)
Student's t :	
central,	
bivariate,	
probabilities	nag_bivariate_students_t (g01hcc)
multivariate,	
probabilities	nag_multi_students_t (g01hdc)
univariate,	
deviates,	
scalar	nag_deviates_students_t (g01fbc)
vectorized	nag_deviates_students_t_vector (g01tbc)
probabilities,	
scalar	nag_prob_students_t (g01ebc)
vectorized	nag_prob_students_t_vector (g01sbc)
non-central,	
probabilities	nag_prob_non_central_students_t (g01gbc)
Studentized range statistic,	
deviates	nag_deviates_studentized_range (g01fmc)
probabilities	nag_prob_studentized_range (g01emc)
von Mises,	
probabilities	nag_prob_von_mises (g01erc)
χ^2 :	
central,	
deviates	nag_deviates_chi_sq (g01fcc)
probabilities	nag_prob_chi_sq (g01ecc)
probability of linear combination	nag_prob_lin_chi_sq (g01jdc)
non-central,	
probabilities	nag_prob_non_central_chi_sq (g01gcc)
probability of linear combination	nag_prob_lin_non_central_chi_sq (g01jcc)
vectorized deviates	nag_deviates_chi_sq_vector (g01tcc)
vectorized probabilities	nag_prob_chi_sq_vector (g01scc)
Scores,	
Normal scores,	
accurate	nag_normal_scores_exact (g01dac)
variance-covariance matrix	nag_normal_scores_var (g01dcc)
Normal scores, ranks or exponential (Savage) scores	nag_ranks_and_scores (g01dhc)

Note: the Student's t , χ^2 , and F functions do not aim to achieve a high degree of accuracy, only about four or five significant figures, but this should be quite sufficient for hypothesis testing. However, both the Student's t and the F -distributions can be transformed to a beta distribution and the χ^2 -distribution can be transformed to a gamma distribution, so a higher accuracy can be obtained by calls to the gamma or beta functions.

Note: `nag_ranks_and_scores` (g01dhc) computes either ranks, approximations to the Normal scores, Normal, or Savage scores for a given sample. `nag_ranks_and_scores` (g01dhc) also gives you control over how it handles tied observations. `nag_normal_scores_exact` (g01dac) computes the Normal scores for a given sample size to a requested accuracy; the scores are returned in ascending order. `nag_normal_scores_exact` (g01dac) can be used if either high accuracy is required or if Normal scores are required for many samples of the same size, in which case you will have to sort the data or scores.

3.1 Working with Streamed or Extremely Large Datasets

The majority of the functions in this chapter are ‘in-core’, that is all the data required must be held in memory prior to calling the function. In some situations this might not be possible, for example, when working with extremely large datasets or where all of the data is not available at once (i.e., the data is being streamed).

There are five functions in this chapter applicable to datasets of this form:

`nag_summary_stats_onevar` (g01atc) computes the mean, variance and the coefficients of skewness and kurtosis for a single variable.

`nag_summary_stats_onevar_combine` (g01auc), takes the results from two calls to `nag_summary_stats_onevar` (g01atc) and combines them, returning the mean, variance and the coefficients of skewness and kurtosis for the combined dataset. This function allows the easy utilization of more than one processor to spread the computational burden inherent in summarising a very large dataset.

`nag_approx_quantiles_fixed` (g01anc) and `nag_approx_quantiles_arbitrary` (g01apc) compute the approximate quantiles for a dataset of known and unknown size respectively.

`nag_moving_average` (g01wac) computes the mean and standard deviation in a rolling window.

In addition, see `nag_sum_sqs` (g02buc) and `nag_sum_sqs_combine` (g02bzc) for functions to summarise two or more variables.

4 Auxiliary Functions Associated with Library Function Arguments

None.

5 Functions Withdrawn or Scheduled for Withdrawal

The following lists all those functions that have been withdrawn since Mark 23 of the Library or are scheduled for withdrawal at one of the next two marks.

Withdrawn Function	Mark of Withdrawal	Replacement Function(s)
<code>nag_summary_stats_1var</code> (g01aac)	26	<code>nag_summary_stats_onevar</code> (g01atc)
<code>nag_deviates_normal_dist</code> (g01cec)	24	<code>nag_deviates_normal</code> (g01fac)

6 References

Hastings N A J and Peacock J B (1975) *Statistical Distributions* Butterworth

Kendall M G and Stuart A (1969) *The Advanced Theory of Statistics (Volume 1)* (3rd Edition) Griffin

Tukey J W (1977) *Exploratory Data Analysis* Addison–Wesley