

# Package: gofedf (via r-universe)

September 6, 2024

**Title** Goodness of Fit Tests Based on Empirical Distribution Functions

**Version** 1.0.0

**Description** Routines that allow the user to run goodness of fit tests based on empirical distribution functions for formal model evaluation in a general likelihood model. In addition, functions are provided to test a sample against Normal or Gamma distributions, validate the normality assumptions in a linear model, and examine the appropriateness of a Gamma distribution in generalized linear models with various link functions.

Michael Arthur Stephens (1976)

<<http://www.jstor.org/stable/2958206>>.

**License** GPL (>= 3)

**URL** <https://github.com/pnickchi/gofedf>

**BugReports** <https://github.com/pnickchi/gofedf/issues>

**Encoding** UTF-8

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.3.1

**Depends** R (>= 2.10)

**Imports** stats, CompQuadForm, MASS, glm2, statmod

**NeedsCompilation** no

**Author** Richard Lockhart [aut], Payman Nickchi [aut, cre]

**Maintainer** Payman Nickchi <[pnickchi@sfu.ca](mailto:pnickchi@sfu.ca)>

**Suggests** knitr, rmarkdown, testthat (>= 3.0.0)

**Config/testthat/edition** 3

**VignetteBuilder** knitr

**Repository** <https://pnickchi.r-universe.dev>

**RemoteUrl** <https://github.com/pnickchi/gofedf>

**RemoteRef** HEAD

**RemoteSha** 36de6ef14c8298cd5ee5e47b5a0548fd5a76374f

## Contents

IGMLE . . . . .	2
IGPIT . . . . .	3
IGScore . . . . .	3
testExponential . . . . .	4
testGamma . . . . .	5
testGLMGamma . . . . .	6
testLMNormal . . . . .	8
testNormal . . . . .	10
testYourModel . . . . .	11
<b>Index</b>	<b>14</b>

---

IGMLE	<i>Compute the maximum likelihood estimate of parameters in Inverse Gaussian distribution with weighted observations.</i>
-------	---

---

### Description

This function is used in [testYourModel](#) function for example purposes.

### Usage

```
IGMLE(obs, ...)
```

### Arguments

obs	a numeric vector of sample observations.
...	a list of additional parameters to define the likelihood.

### Value

The function compute the MLE of parameters in Inverse Gaussian distribution and returns a vector of estimates. The first and second elements of the vector are MLE of the mean and shape, respectively.

---

IGPIT	<i>Compute the probability integral transformed values for a sample from Inverse Gaussian distribution.</i>
-------	---

---

**Description**

This function is used in [testYourModel](#) function for example purposes.

**Usage**

```
IGPIT(obs, ...)
```

**Arguments**

obs	A numeric vector of sample observations.
...	A list of additional parameters to define the likelihood.

**Value**

A numeric vector of probability integral transformed values of sample observations.

---

IGScore	<i>Compute the score function of the Inverse Gaussian distribution based on a sample.</i>
---------	---

---

**Description**

This function is used in [testYourModel](#) function for example purposes.

**Usage**

```
IGScore(obs, ...)
```

**Arguments**

obs	a numeric vector of sample observations.
...	a list of additional parameters to define the likelihood.

**Value**

The score matrix with n rows (number of sample observations) and 2 columns (mean and shape).

---

testExponential      *Apply Goodness of Fit Test for Exponential Distribution*

---

### Description

Performs the goodness-of-fit test based on empirical distribution function to check if an i.i.d sample follows an Exponential distribution.

### Usage

```
testExponential(
  x,
  discretize = FALSE,
  ngrid = length(x),
  gridpit = FALSE,
  hessian = FALSE,
  method = "cvm"
)
```

### Arguments

x	a non-empty numeric vector of sample data.
discretize	If TRUE, the covariance function of $W_n(u)$ process is evaluated at some data points (see ngrid and gridpit), and the integral equation is replaced by a matrix equation. If FALSE (the default value), the covariance function is first estimated, and then the integral equation is solved to find the eigenvalues. The results of our simulations recommend using the estimated covariance for solving the integral equation. The parameters ngrid, gridpit, and hessian are only relevant when discretize = TRUE.
ngrid	the number of equally spaced points to discretize the (0,1) interval for computing the covariance function.
gridpit	logical. If TRUE (the default value), the parameter ngrid is ignored and (0,1) interval is divided based on probability integral transformed values obtained from the sample. If FALSE, the interval is divided into ngrid equally spaced points for computing the covariance function.
hessian	logical. If TRUE the Fisher information matrix is estimated by the observed Hessian Matrix based on the sample. If FALSE (the default value) the Fisher information matrix is estimated by the variance of the observed score matrix.
method	a character string indicating which goodness-of-fit statistic is to be computed. The default value is 'cvm' for the Cramer-von-Mises statistic. Other options include 'ad' for the Anderson-Darling statistic, and 'both' to compute both cvm and ad.

**Value**

A list of two containing the following components:

- **Statistic:** the value of goodness-of-fit statistic.
- **p-value:** the approximate p-value for the goodness-of-fit test. if `method = 'cvm'` or `method = 'ad'`, it returns a numeric value for the statistic and p-value. If `method = 'both'`, it returns a numeric vector with two elements and one for each statistic.

**Examples**

```
set.seed(123)
n <- 50
sim_data <- rexp(n, rate = 2)
testExponential(x = sim_data)
```

---

testGamma

*Apply Goodness of Fit Test for Gamma Distribution*


---

**Description**

Performs the goodness-of-fit test based on empirical distribution function to check if an i.i.d sample follows a Gamma distribution.

**Usage**

```
testGamma(
  x,
  discretize = FALSE,
  ngrid = length(x),
  gridpit = FALSE,
  hessian = FALSE,
  rate = TRUE,
  method = "cvm"
)
```

**Arguments**

<code>x</code>	a non-empty numeric vector of sample data.
<code>discretize</code>	If TRUE, the covariance function of $W_n(u)$ process is evaluated at some data points (see <code>ngrid</code> and <code>gridpit</code> ), and the integral equation is replaced by a matrix equation. If FALSE (the default value), the covariance function is first estimated, and then the integral equation is solved to find the eigenvalues. The results of our simulations recommend using the estimated covariance for solving the integral equation. The parameters <code>ngrid</code> , <code>gridpit</code> , and <code>hessian</code> are only relevant when <code>discretize = TRUE</code> .
<code>ngrid</code>	the number of equally spaced points to discretize the (0,1) interval for computing the covariance function.

gridpit	logical. If TRUE (the default value), the parameter ngrid is ignored and (0,1) interval is divided based on probability integral transformed values obtained from the sample. If FALSE, the interval is divided into ngrid equally spaced points for computing the covariance function.
hessian	logical. If TRUE the Fisher information matrix is estimated by the observed Hessian Matrix based on the sample. If FALSE (the default value) the Fisher information matrix is estimated by the variance of the observed score matrix.
rate	logical. If TRUE (the default value), the rate is estimated in Gamma distribution. If FALSE, scale is estimated. See <a href="#">GammaDist</a> for more details.
method	a character string indicating which goodness-of-fit statistic is to be computed. The default value is 'cvm' for the Cramer-von-Mises statistic. Other options include 'ad' for the Anderson-Darling statistic, and 'both' to compute both cvm and ad.

### Value

A list of two containing the following components:

- Statistic: the value of goodness-of-fit statistic.
- p-value: the approximate p-value for the goodness-of-fit test. if method = 'cvm' or method = 'ad', it returns a numeric value for the statistic and p-value. If method = 'both', it returns a numeric vector with two elements and one for each statistic.

### Examples

```
set.seed(123)
sim_data <- rgamma(n = 50, shape = 3)
testGamma(x = sim_data)
sim_data <- runif(n = 50)
testGamma(x = sim_data)
```

---

testGLMGamma	<i>Apply Goodness of Fit Test to the Residuals of a Generalized Linear Model with Gamma Link Function</i>
--------------	---

---

### Description

testGLMGamma is used to check the validity of Gamma assumption for the response variable when fitting generalized linear model. Common link functions in [glm](#) can be used here.

### Usage

```
testGLMGamma(
  x,
  y,
  fit = NULL,
  l = "log",
```

```

    discretize = FALSE,
    ngrid = length(y),
    gridpit = TRUE,
    hessian = FALSE,
    start.value = NULL,
    control = NULL,
    method = "cvm"
)

```

### Arguments

<code>x</code>	is either a numeric vector or a design matrix. In the design matrix, rows indicate observations and columns presents covariats.
<code>y</code>	is a vector of numeric values with the same number of observations or number of rows as <code>x</code> .
<code>fit</code>	is an object of class <code>glm</code> and its default value is <code>NULL</code> . If a fit of class <code>glm</code> is provided, the arguments <code>x</code> , <code>y</code> , and <code>l</code> will be ignored. We recommend using <code>glm2</code> function from <code>glm2</code> package since it provides better convergence while optimizing the likelihood to estimate coefficients of the model by IWLS method. It is required to return design matrix by <code>x = TRUE</code> in <code>glm</code> or <code>glm2</code> function. For more information on how to do this, refer to the help documentation for the <code>glm</code> or <code>glm2</code> function.
<code>l</code>	a character vector indicating the link function that should be used for Gamma family. Acceptable link functions for Gamma family are inverse, identity and log. For more details see <code>Gamma</code> from <code>stats</code> package.
<code>discretize</code>	If <code>TRUE</code> , the covariance function of $W_n(u)$ process is evaluated at some data points (see <code>ngrid</code> and <code>gridpit</code> ), and the integral equation is replaced by a matrix equation. If <code>FALSE</code> (the default value), the covariance function is first estimated, and then the integral equation is solved to find the eigenvalues. The results of our simulations recommend using the estimated covariance for solving the integral equation. The parameters <code>ngrid</code> , <code>gridpit</code> , and <code>hessian</code> are only relevant when <code>discretize = TRUE</code> .
<code>ngrid</code>	the number of equally spaced points to discretize the (0,1) interval for computing the covariance function.
<code>gridpit</code>	logical. If <code>TRUE</code> (the default value), the parameter <code>ngrid</code> is ignored and (0,1) interval is divided based on probability integral transformed values obtained from the sample. If <code>FALSE</code> , the interval is divided into <code>ngrid</code> equally spaced points for computing the covariance function.
<code>hessian</code>	logical. If <code>TRUE</code> the Fisher information matrix is estimated by the observed Hessian Matrix based on the sample. If <code>FALSE</code> (the default value) the Fisher information matrix is estimated by the variance of the observed score matrix.
<code>start.value</code>	a numeric value or vector. This is the same as <code>start</code> argument in <code>glm</code> or <code>glm2</code> . The value is a starting point in iteratively reweighted least squares (IRLS) algorithm for estimating the MLE of coefficients in the model.
<code>control</code>	a list of parameters to control the fitting process in <code>glm</code> or <code>glm2</code> function. For more details, see <code>glm.control</code> .

**method** a character string indicating which goodness-of-fit statistic is to be computed. The default value is 'cvm' for the Cramer-von-Mises statistic. Other options include 'ad' for the Anderson-Darling statistic, and 'both' to compute both cvm and ad.

### Value

A list of two containing the following components:

- **Statistic:** the value of goodness-of-fit statistic.
- **p-value:** the approximate p-value for the goodness-of-fit test. if `method = 'cvm'` or `method = 'ad'`, it returns a numeric value for the statistic and p-value. If `method = 'both'`, it returns a numeric vector with two elements and one for each statistic.
- **converged:** logical to indicate if the IWLS algorithm have converged or not.

### Examples

```
set.seed(123)
n <- 50
p <- 5
x <- matrix( rnorm(n*p, mean = 10, sd = 0.1), nrow = n, ncol = p)
b <- runif(p)
e <- rgamma(n, shape = 3)
y <- exp(x %*% b) * e
testGLMGamma(x, y, l = 'log')
myfit <- glm(y ~ x, family = Gamma('log'), x = TRUE, y = TRUE)
testGLMGamma(fit = myfit)
```

---

testLMNormal

*Apply Goodness of Fit Test to Residuals of a Linear Model*

---

### Description

testLMNormal is used to check the normality assumption of residuals in a linear model. This function can take the response variable and design matrix, fit a linear model, and apply the goodness-of-fit test. Conveniently, it can take an object of class "lm" and directly applies the goodness-of-fit test. The function returns a goodness-of-fit statistic along with an approximate p-value.

### Usage

```
testLMNormal(
  x,
  y,
  fit = NULL,
  discretize = FALSE,
  ngrid = length(y),
  gridpit = TRUE,
  hessian = FALSE,
  method = "cvm"
)
```

**Arguments**

<code>x</code>	is either a numeric vector or a design matrix. In the design matrix, rows indicate observations and columns presents covariates.
<code>y</code>	is a vector of numeric values with the same number of observations or number of rows as <code>x</code> .
<code>fit</code>	an object of class "lm" returned by <code>lm</code> function in <code>stats</code> package. The default value of <code>fit</code> is <code>NULL</code> . If any object is provided, <code>x</code> and <code>y</code> will be ignored and the class of object is checked. If you pass an object to <code>fit</code> make sure to return the design matrix by setting <code>x = TRUE</code> and the response variable by setting <code>y = TRUE</code> in <code>lm</code> function. To read more about this see the help documentation for <code>lm</code> function or see the example below.
<code>discretize</code>	If <code>TRUE</code> , the covariance function of $W_n(u)$ process is evaluated at some data points (see <code>ngrid</code> and <code>gridpit</code> ), and the integral equation is replaced by a matrix equation. If <code>FALSE</code> (the default value), the covariance function is first estimated, and then the integral equation is solved to find the eigenvalues. The results of our simulations recommend using the estimated covariance for solving the integral equation. The parameters <code>ngrid</code> , <code>gridpit</code> , and <code>hessian</code> are only relevant when <code>discretize = TRUE</code> .
<code>ngrid</code>	the number of equally spaced points to discretize the (0,1) interval for computing the covariance function.
<code>gridpit</code>	logical. If <code>TRUE</code> (the default value), the parameter <code>ngrid</code> is ignored and (0,1) interval is divided based on probability integral transformed values obtained from the sample. If <code>FALSE</code> , the interval is divided into <code>ngrid</code> equally spaced points for computing the covariance function.
<code>hessian</code>	logical. If <code>TRUE</code> the Fisher information matrix is estimated by the observed Hessian Matrix based on the sample. If <code>FALSE</code> (the default value) the Fisher information matrix is estimated by the variance of the observed score matrix.
<code>method</code>	a character string indicating which goodness-of-fit statistic is to be computed. The default value is <code>'cvm'</code> for the Cramer-von-Mises statistic. Other options include <code>'ad'</code> for the Anderson-Darling statistic, and <code>'both'</code> to compute both <code>cvm</code> and <code>ad</code> .

**Value**

A list of two containing the following components:

- **Statistic:** the value of goodness-of-fit statistic.
- **p-value:** the approximate p-value for the goodness-of-fit test. if `method = 'cvm'` or `method = 'ad'`, it returns a numeric value for the statistic and p-value. If `method = 'both'`, it returns a numeric vector with two elements and one for each statistic.

**Examples**

```
set.seed(123)
n <- 50
p <- 5
x <- matrix( runif(n*p), nrow = n, ncol = p)
```

```

e <- rnorm(n)
b <- runif(p)
y <- x %*% b + e
testLMNormal(x, y)
# Or pass lm.fit object directly:
lm.fit <- lm(y ~ x, x = TRUE, y = TRUE)
testLMNormal(fit = lm.fit)

```

---

testNormal

*Apply Goodness of Fit Test for Normal Distribution*


---

### Description

Performs the goodness-of-fit test based on empirical distribution function to check if an i.i.d sample follows a Normal distribution.

### Usage

```

testNormal(
  x,
  discretize = FALSE,
  ngrid = length(x),
  gridpit = TRUE,
  hessian = FALSE,
  method = "cvm"
)

```

### Arguments

x	a non-empty numeric vector of sample data.
discretize	If TRUE, the covariance function of $W_n(u)$ process is evaluated at some data points (see ngrid and gridpit), and the integral equation is replaced by a matrix equation. If FALSE (the default value), the covariance function is first estimated, and then the integral equation is solved to find the eigenvalues. The results of our simulations recommend using the estimated covariance for solving the integral equation. The parameters ngrid, gridpit, and hessian are only relevant when discretize = TRUE.
ngrid	the number of equally spaced points to discretize the (0,1) interval for computing the covariance function.
gridpit	logical. If TRUE (the default value), the parameter ngrid is ignored and (0,1) interval is divided based on probability inverse transformed values obtained from the sample. If FALSE, the interval is divided into ngrid equally spaced points for computing the covariance function.
hessian	logical. If TRUE the Fisher information matrix is estimated by the observed Hessian Matrix based on the sample. If FALSE (the default value) the Fisher information matrix is estimated by the variance of the observed score matrix.

**method** a character string indicating which goodness-of-fit statistic is to be computed. The default value is 'cvm' for the Cramer-von-Mises statistic. Other options include 'ad' for the Anderson-Darling statistic, and 'both' to compute both cvm and ad.

### Value

A list of two containing the following components:

- **Statistic:** the value of goodness-of-fit statistic.
- **p-value:** the approximate p-value for the goodness-of-fit test. if method = 'cvm' or method = 'ad', it returns a numeric value for the statistic and p-value. If method = 'both', it returns a numeric vector with two elements and one for each statistic.

### Examples

```
set.seed(123)
sim_data <- rnorm(n = 50)
testNormal(x = sim_data)
sim_data <- rgamma(n = 50, shape = 3)
testNormal(x = sim_data)
```

---

testYourModel

*Apply the Goodness of Fit Test Based on Empirical Distribution Function to Any Likelihood Model.*

---

### Description

This function applies the goodness-of-fit test based on empirical distribution function. It requires certain inputs depending on whether the model involves parameter estimation or not. If the model is known and there is no parameter estimation, the function requires the probability transformed (or pit) values of the sample. This ought to be a numeric vector. If there is parameter estimation in the model, the function additionally requires the score as a matrix with n rows and p columns, where n is the sample size and p is the number of estimated parameters. The function checks if the sum of columns in score is near zero at the estimated parameter (which is assumed to be the maximum likelihood estimate).

### Usage

```
testYourModel(
  pit,
  score = NULL,
  discretize = FALSE,
  ngrid = length(pit),
  gridpit = TRUE,
  precision = 1e-09,
  method = "cvm"
)
```

**Arguments**

pit	The probability transformed (or pit) values of the sample which ought to be a numeric vector.
score	The default value is null and refers to no parameter estimation case. If there is parameter estimation, the score must be a matrix with n rows and p columns, where n is the sample size and p is the number of estimated parameters.
discretize	If TRUE, the covariance function of $W_n(u)$ process is evaluated at some data points (see ngrid and gridpit), and the integral equation is replaced by a matrix equation. If FALSE (the default value), the covariance function is first estimated, and then the integral equation is solved to find the eigenvalues. The results of our simulations recommend using the estimated covariance for solving the integral equation. The parameters ngrid, gridpit, and hessian are only relevant when discretize = TRUE.
ngrid	The number of equally spaced points to discretize the (0,1)interval for computing the covariance function.
gridpit	logical. If TRUE (the default value), the parameter ngrid is ignored and (0,1) interval is divided based on probability integral transformed values obtained from the sample. If FALSE, the interval is divided into ngrid equally spaced points for computing the covariance function.
precision	The theory behind goodness-of-fit test based on empirical distribution function (edf) works well if the MLE is indeed the root of derivative of log likelihood function. A precision of 1e-9 (default value) is used to check this. A warning message is generated if the score evaluated at MLE is not close enough to zero.
method	a character string indicating which goodness-of-fit statistic is to be computed. The default value is 'cvm' for the Cramer-von-Mises statistic. Other options include 'ad' for the Anderson-Darling statistic, and 'both' to compute both cvm and ad.

**Value**

A list of two containing the following components:

- Statistic: the value of goodness-of-fit statistic.
- p-value: the approximate p-value for the goodness-of-fit test. if method = 'cvm' or method = 'ad', it returns a numeric value for the statistic and p-value. If method = 'both', it returns a numeric vector with two elements and one for each statistic.

**Examples**

```
# Example: Inverse Gaussian (IG) distribution with weights

# Set the seed to reproduce example.
set.seed(123)

# Set the sample size
n <- 50
```

```
# Assign weights
weights <- rep(1.5,n)

# Set mean and shape parameters for IG distribution.
mio      <- 2
lambda  <- 2

# Generate a random sample from IG distribution with weighted shape.
sim_data <- statmod::rinvgauss(n, mean = mio, shape = lambda * weights)

# Compute MLE of parameters, score matrix, and pit values.
theta_hat <- IGMLE(obs = sim_data, w = weights)
ScoreMatrix <- IGScore(obs = sim_data, w = weights, mle = theta_hat)
pitvalues <- IGPIT(obs = sim_data, w = weights, mle = theta_hat)

# Apply the goodness-of-fit test.
testYourModel(pit = pitvalues, score = ScoreMatrix)
```

# Index

Gamma, [7](#)  
GammaDist, [6](#)  
glm, [6](#), [7](#)  
glm.control, [7](#)  
glm2, [7](#)

IGMLE, [2](#)  
IGPIT, [3](#)  
IGScore, [3](#)

lm, [9](#)

stats, [9](#)

testExponential, [4](#)  
testGamma, [5](#)  
testGLMGamma, [6](#)  
testLMNormal, [8](#)  
testNormal, [10](#)  
testYourModel, [2](#), [3](#), [11](#)