

# Package ‘ganGenerativeData’

March 25, 2023

**Type** Package

**Title** Generate Generative Data for a Data Source

**Version** 1.4.3

**Date** 2023-03-25

**Author** Werner Mueller

**Maintainer** Werner Mueller <werner.mueller5@chello.at>

## Description

Generative Adversarial Networks are applied to generate generative data for a data source. In iterative training steps the distribution of generated data converges to that of the data source. Direct applications of generative data are the created functions for data classifying and missing data completion. Reference: Goodfellow et al. (2014) <[arXiv:1406.2661v1](#)>.

**License** GPL (>= 2)

**Imports** Rcpp (>= 1.0.3), tensorflow (>= 2.0.0)

**LinkingTo** Rcpp

**RoxygenNote** 7.2.3

**SystemRequirements** TensorFlow (<https://www.tensorflow.org>)

**NeedsCompilation** yes

**Encoding** UTF-8

**Repository** CRAN

**Date/Publication** 2023-03-25 18:00:02 UTC

## R topics documented:

ganGenerativeData-package	2
dsActivateColumns	9
dsCreateWithDataFrame	9
dsDeactivateColumns	10
dsGetActiveColumnNames	10
dsGetInactiveColumnNames	11
dsGetNumberOfRows	11
dsGetRow	12

dsRead . . . . .	12
dsWrite . . . . .	13
gdCalculateDensityValue . . . . .	14
gdCalculateDensityValueQuantile . . . . .	14
gdCalculateDensityValues . . . . .	15
gdComplete . . . . .	16
gdGenerate . . . . .	16
gdGenerateParameters . . . . .	17
gdGetNumberOfRows . . . . .	18
gdGetRow . . . . .	19
gdKNearestNeighbors . . . . .	19
gdPlotDataSourceParameters . . . . .	20
gdPlotParameters . . . . .	21
gdPlotProjection . . . . .	21
gdRead . . . . .	22
gdWriteSubset . . . . .	23

<b>Index</b>	<b>24</b>
--------------	-----------

---

ganGenerativeData-package

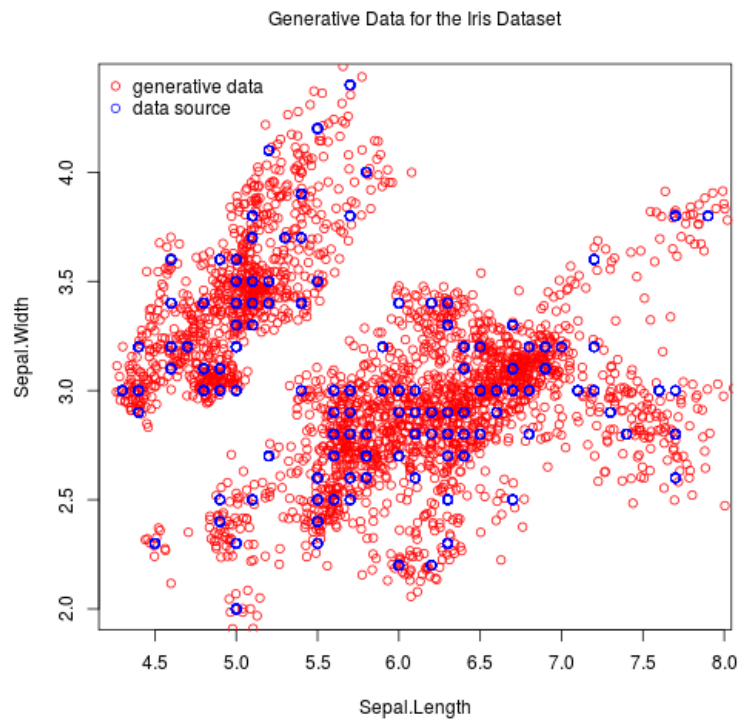
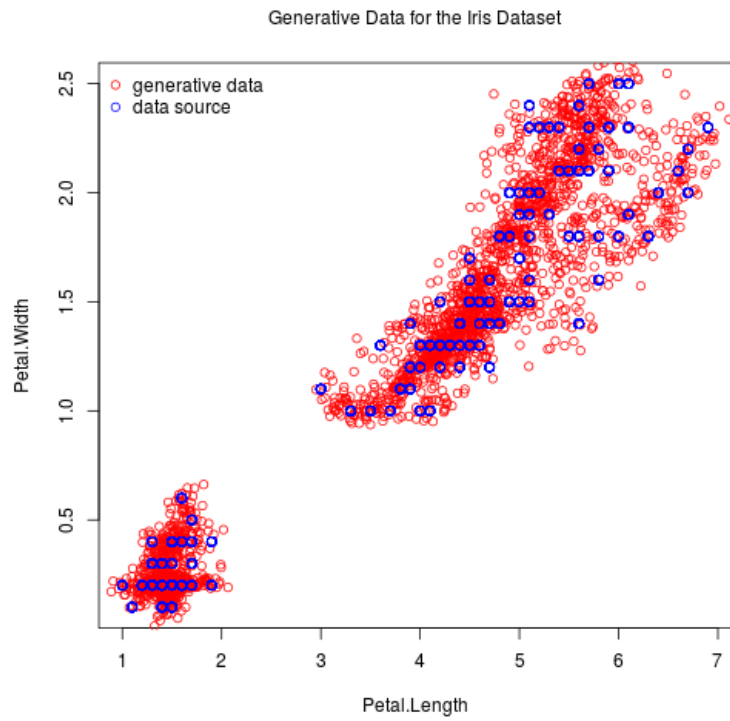
*Generate generative data for a data source*

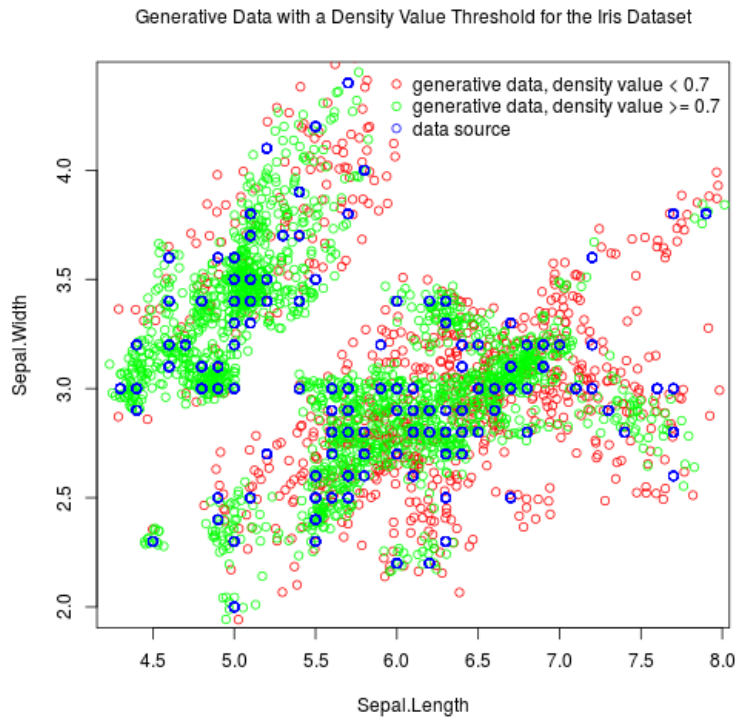
---

## Description

Generative Adversarial Networks are applied to generate generative data for a data source. In iterative training steps the distribution of generated data converges to that of the data source. Direct applications of generative data are the created functions for data classifying and missing data completion.

The inserted images show two-dimensional projections of generative data for the iris dataset:





## Details

The API includes functions for topics "definition of data source" and "generation of generative data". Main function of first topic is `dsCreateWithDataFrame()` which creates a data source with passed data frame. Main function of second topic is `gdGenerate()` which generates generative data for a data source.

### 1. Definition of data source

`dsCreateWithDataFrame()` Create a data source with passed data frame.

`dsActivateColumns()` Activate columns of a data source in order to include them in generation of generative data. By default columns are active.

`dsDeactivateColumns()` Deactivate columns of a data source in order to exclude them in generation of generative data. Note that in this version only columns of type R-class numeric and R-type double can be used in generation of generative data. All columns of other type have to be deactivated.

`dsGetActiveColumnNames()` Get names of active columns of a data source.

`dsGetInactiveColumnNames()` Get names of inactive columns of a data source.

`dsWrite()` Write created data source including settings of active columns to a file in binary format. This file will be used as input in functions of topic "generation of generative data".

`dsRead()` Read a data source from a file that was written with `dsWrite()`.

`dsGetNumberOfRows()` Get number of rows in a data source.

`dsGetRow()` Get a row in a data source.

### 2. Generation of generative data

`gdGenerateParameters()` Specify parameters for generation of generative data.

`gdGenerate()` Read a data source from a file, generate generative data for the data source in iterative training steps and write generated data to a file in binary format.

`gdCalculateDensityValues()` Read generative data from a file, calculate density values and write generative data with density values to original file.

`gdRead()` Read generative data and data source from specified files.

gdPlotParameters() Specify plot parameters for generative data.

gdPlotDataSourceParameters() Specify plot parameters for data source.

gdPlotProjection() Create an image file containing two-dimensional projections of generative data and data source.

gdGetNumberOfRows() Get number of rows in generative data.

gdGetRow() Get a row in generative data.

gdCalculateDensityValue() Calculate density value for a data record.

gdCalculateDensityValueQuantile() Calculate density value quantile for a percent value.

gdKNearestNeighbors() Search for k nearest neighbors in generative data.

gdComplete() Complete incomplete data record.

gdWriteSubset() Write subset of generative data.

### Author(s)

Werner Mueller  
Maintainer: Werner Mueller <werner.mueller5@chello.at>

### References

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio (2014), "*Generative Adversarial Nets*", <arXiv:1406.2661v1>

### Examples

```
# Environment used for execution of examples:

# Operating system: Ubuntu 22.04.1
# Compiler: g++ 11.3.0 (supports C++17 standard)
# R applications: R 4.1.2, RStudio 2022.02.2
# Installed packages: 'Rcpp' 1.0.10, 'tensorflow' 2.11.0,
# 'ganGenerativeData' 1.4.2

# Package 'tensorflow' provides an interface to machine learning framework
# TensorFlow. To complete the installation function install_tensorflow() has to
# be called.
## Not run:
library(tensorflow)
install_tensorflow()
```

```
## End(Not run)

# Generate generative data for the iris dataset

# Load library
library(ganGenerativeData)

# 1. Definition of data source for the iris dataset

# Create a data source with iris data frame.
dsCreateWithDataFrame(iris)

# Deactivate the column with name Species and index 5 in order to exclude it in
# generation of generative data.
dsDeactivateColumns(c(5))

# Get the active column names: Sepal.Length, Sepal.Width, Petal.Length,
# Petal.Width.
dsGetActiveColumnNames()

# Write the data source including settings of active columns to file
# "iris4d.bin" in binary format.
## Not run:
dsWrite("ds.bin")
## End(Not run)

# 2. Generation of generative data for the iris data source

# Read data source from file "ds.bin", generate generative data in iterative
# training steps (in tests 50000 iterations are used) and write generated
# generative data to file "gd.bin".
## Not run:
gdGenerate("gd.bin", "ds.bin", c(1, 2), gdGenerateParameters(2500))
## End(Not run)

# Read generative data from file "gd.bin", calculate density values and
# write generative data with density values to original file.
## Not run:
gdCalculateDensityValues("gd.bin")
## End(Not run)

# Read generative data from file "gd.bin" and data source from "ds.bin"
## Not run:
gdRead("gd.bin", "ds.bin")
## End(Not run)

# Create an image showing two-dimensional projections of generative data and
# data source for column indices 3, 4 and write it to file "gd34d.png"
## Not run:
gdPlotProjection("gd34d.png",
"Generative Data for the Iris Dataset",
c(3, 4),
gdPlotParameters(25),
```

```
gdPlotDataSourceParameters(100))
## End(Not run)

# Create an image showing two-dimensional projections of generative data and
# data source for column indices 3, 4 with density value threshold 0.71 and
# write it to file "gd34ddv.png"
## Not run:
gdPlotProjection("gd34ddv.png",
  "Generative Data with a Density Value Threshold for the Iris Dataset",
  c(3, 4),
  gdPlotParameters(25, c(0.71), c("red", "green")),
  gdPlotDataSourceParameters(100))
## End(Not run)

# Get number of rows in generative data
## Not run:
gdGetNumberOfRows()
## End(Not run)

# Get row with index 1000 in generative data
## Not run:
gdGetRow(1000)
## End(Not run)

# Calculate density value for a data record
## Not run:
gdCalculateDensityValue(list(6.1, 2.6, 5.6, 1.4))
## End(Not run)

# Calculate density value quantile for 50 percent
## Not run:
gdCalculateDensityValueQuantile(50)
## End(Not run)

# Search for k nearest neighbors for a data record
## Not run:
gdKNearestNeighbors(list(5.1, 3.5, 1.4, 0.2), 3)
## End(Not run)

# Complete incomplete data record containing an NA value
## Not run:
gdComplete(list(5.1, 3.5, 1.4, NA))
## End(Not run)

# Write subset containing 50 percent of randomly selected rows of
# generative data
## Not run:
gdRead("gd.bin")
gdWriteSubset("gds.bin", 50)
## End(Not run)
```



---

dsActivateColumns	<i>Activate columns</i>
-------------------	-------------------------

---

**Description**

Activate columns of a data source in order to include them in generation of generative data. By default columns are active.

**Usage**

```
dsActivateColumns(columnVector)
```

**Arguments**

columnVector    Vector of column indices

**Value**

None

**Examples**

```
dsCreateWithDataFrame(iris)
dsGetActiveColumnNames()
dsDeactivateColumns(c(5))
dsGetActiveColumnNames()
dsActivateColumns(c(5))
dsGetActiveColumnNames()
```

---

dsCreateWithDataFrame	<i>Create a data source with passed data frame</i>
-----------------------	--

---

**Description**

Create a data source with passed data frame.

**Usage**

```
dsCreateWithDataFrame(dataFrame)
```

**Arguments**

dataFrame        Name of data frame

**Value**

None

**Examples**

```
dsCreateWithDataFrame(iris)
```

---

```
dsDeactivateColumns    Deactivate columns
```

---

**Description**

Deactivate columns of a data source in order to exclude them in generation of generative data. Note that in this version only columns of type R-class numeric and R-type double can be used in generator of generative data. All columns of other type have to be deactivated.

**Usage**

```
dsDeactivateColumns(columnVector)
```

**Arguments**

columnVector    Vector of column indices

**Value**

None

**Examples**

```
dsCreateWithDataFrame(iris)
dsDeactivateColumns(c(5))
dsGetInactiveColumnNames()
```

---

```
dsGetActiveColumnNames
```

*Get active column names*

---

**Description**

Get active column names of a data source

**Usage**

```
dsGetActiveColumnNames()
```

**Value**

Vector of names of active columns

**Examples**

```
dsCreateWithDataFrame(iris)
dsDeactivateColumns(c(5))
dsGetActiveColumnNames()
```

---

```
dsGetInactiveColumnNames
    Get inactive column names
```

---

**Description**

Get inactive column names of a data source

**Usage**

```
dsGetInactiveColumnNames()
```

**Value**

Vector of names of inactive columns

**Examples**

```
dsCreateWithDataFrame(iris)
dsDeactivateColumns(c(5))
dsGetInactiveColumnNames()
```

---

```
dsGetNumberOfRows    Get number of rows
```

---

**Description**

Get number of rows in a data source

**Usage**

```
dsGetNumberOfRows()
```

**Value**

Number of rows

**Examples**

```
dsCreateWithDataFrame(iris)
dsGetNumberOfRows()
```

---

dsGetRow	<i>Get a row in a data source</i>
----------	-----------------------------------

---

**Description**

Get a row in a data source for a row index.

**Usage**

```
dsGetRow(index)
```

**Arguments**

index	Index of row
-------	--------------

**Value**

List containing row in data source

**Examples**

```
dsCreateWithDataFrame(iris)  
dsGetRow(1)
```

---

dsRead	<i>Read a data source from file</i>
--------	-------------------------------------

---

**Description**

Read a data source from a file in binary format

**Usage**

```
dsRead(fileName)
```

**Arguments**

fileName	Name of data source file
----------	--------------------------

**Value**

None

**Examples**

```
## Not run:  
dsCreateWithDataFrame(iris)  
dsDeactivateColumns(c(5))  
dsWrite("ds.bin")  
dsRead("ds.bin")  
## End(Not run)
```

---

dsWrite

*Write a data source to file*

---

**Description**

Write a data source including settings of active columns to a file in binary format. This file will be used as input in functions for generation of generative data.

**Usage**

```
dsWrite(fileName)
```

**Arguments**

fileName      Name of data source file

**Value**

None

**Examples**

```
## Not run:  
dsCreateWithDataFrame(iris)  
dsDeactivateColumns(c(5))  
dsWrite("ds.bin")  
## End(Not run)
```

gdCalculateDensityValue

*Calculate density value for a data record*

---

### Description

Calculate density value for a data record. By default for the calculation a linear search is performed on generative data. When a search tree is used search is performed on a tree for generative data which is built once in the first function call.

### Usage

```
gdCalculateDensityValue(dataRecord, useSearchTree = FALSE)
```

### Arguments

dataRecord      List containing a data record  
useSearchTree    Boolean value indicating if a search tree should be used.

### Value

Normalized density value number

### Examples

```
## Not run:  
gdRead("gd.bin")  
gdCalculateDensityValue(list(6.1, 2.6, 5.6, 1.4))  
## End(Not run)
```

---

gdCalculateDensityValueQuantile

*Calculate density value quantile*

---

### Description

Calculate density value quantile for a percent value.

### Usage

```
gdCalculateDensityValueQuantile(percent)
```

### Arguments

percent          Percent value

**Value**

Normalized density value quantile number

**Examples**

```
## Not run:  
gdRead("gd.bin")  
gdCalculateDensityValueQuantile(50)  
## End(Not run)
```

---

gdCalculateDensityValues

*Calculate density values for generative data*

---

**Description**

Read generative data from a file, calculate density values and write generative data with density values to original file. Calculated density values are used to classify generative data. In function gdPlotParameters() density value thresholds with assigned colors can be passed to draw generative data for different density value ranges.

**Usage**

```
gdCalculateDensityValues(generativeDataFileName)
```

**Arguments**

generativeDataFileName  
Name of generative data file name

**Value**

None

**Examples**

```
## Not run:  
gdCalculateDensityValues("gd.bin")  
## End(Not run)
```

---

gdComplete	<i>Complete incomplete data record</i>
------------	--

---

**Description**

Search for first nearest neighbor in generative data for incomplete data record containing NA values. Found row in generative data is then used to replace NA values in incomplete data record. This function calls gdKNearestNeighbors() with parameter k equal to 1.

**Usage**

```
gdComplete(dataRecord, useSearchTree = FALSE)
```

**Arguments**

dataRecord	List containing incomplete data record
useSearchTree	Boolean value indicating if a search tree should be used.

**Value**

List containing completed data record

**Examples**

```
## Not run:
gdRead("gd.bin")
gdComplete(list(5.1, 3.5, 1.4, NA))
## End(Not run)
```

---

gdGenerate	<i>Generate generative data for a data source</i>
------------	---

---

**Description**

Read a data source from a file, generate generative data for the data source in iterative training steps and write generated data to a file in binary format. When a higher number of iterations is used the distribution of generated data gets closer to that of the data source.

**Usage**

```
gdGenerate(
  generativeDataFileName,
  dataSourceFileName,
  columnIndices,
  generateParameters = gdGenerateParameters(numberOfIterations = 25000,
  numberOfHiddenLayerUnits = 1024, learningRate = 1e-04, keepProbability = 0.95,
  collectBeginningAtIteration = 1)
)
```



**Arguments**

generativeDataFileName  
Name of generative data file

dataSourceFileName  
Name of data source file

columnIndices  
Vector of two column indices that are used to plot two-dimensional projections of normalized generated generative data and data source for a training step. Indices refer to indices of active columns of data source. Plotting can be disabled by passing NULL or an empty vector.

generateParameters  
Generation of generative data parameters, see function gdGenerateParameters().

**Value**

None

**Examples**

```
## Not run:
generateParameters <- gdGenerateParameters(numberOfIterations = 5000)
gdGenerate("gd.bin", "ds.bin", c(1, 2), generateParameters)
## End(Not run)
```

---

gdGenerateParameters *Specify parameters for generation of generative data*

---

**Description**

Specify parameters for training of neural networks used for generation of generative data. These parameters are passed to function gdGenerate().

**Usage**

```
gdGenerateParameters(
  numberOfIterations = 50000,
  numberOfHiddenLayerUnits = 1024,
  learningRate = 1e-04,
  keepProbability = 0.95,
  collectBeginningAtIteration = 1
)
```

**Arguments**

numberOfIterations  
Number of training steps

numberOfHiddenLayerUnits  
Number of hidden layer units

`learningRate`    Learning rate for training of neural networks  
`keepProbability`    Value in the range of 0 to 1 which is used to train generalized neural networks.  
`collectBeginningAtIteration`    Collect generative data beginning at iteration

**Value**

List of parameters for generation of generative data

**Examples**

```
## Not run:  
generateParameters <- gdGenerateParameters(numberOfIterations = 5000)  
## End(Not run)
```

---

`gdGetNumberOfRows`    *Get number of rows*

---

**Description**

Get number of rows in generative data

**Usage**

```
gdGetNumberOfRows()
```

**Value**

Number of rows

**Examples**

```
## Not run:  
gdRead("gd.bin")  
gdGetNumberOfRows()  
## End(Not run)
```

---

gdGetRow	<i>Get a row in generative data</i>
----------	-------------------------------------

---

**Description**

Get a row in generative data for a row index

**Usage**

```
gdGetRow(index)
```

**Arguments**

index	Index of row
-------	--------------

**Value**

List containing row in generative data

**Examples**

```
## Not run:  
gdRead("gd.bin")  
gdGetRow(1000)  
## End(Not run)
```

---

gdKNearestNeighbors	<i>Search for k nearest neighbors</i>
---------------------	---------------------------------------

---

**Description**

Search for k nearest neighbors in generative data for a data record. When the data record contains NA values only the non-NA values are considered in search. By default a linear search is performed. When a search tree is used search is performed on a tree which is built once in the first function call. Building a tree is also triggered when NA values in data records change in subsequent function calls.

**Usage**

```
gdKNearestNeighbors(dataRecord, k = 1L, useSearchTree = FALSE)
```

**Arguments**

dataRecord	List containing a data record
k	Number of nearest neighbors
useSearchTree	Boolean value indicating if a search tree should be used.

**Value**

A list of rows in generative data

**Examples**

```
## Not run:  
gdRead("gd.bin")  
gdKNearestNeighbors(list(5.1, 3.5, 1.4, 0.2), 3)  
## End(Not run)
```

---

gdPlotDataSourceParameters

*Specify plot parameters for data source*

---

**Description**

Specify plot parameters for data source passed to function gdPlotProjection().

**Usage**

```
gdPlotDataSourceParameters(percent = 100, color = "blue")
```

**Arguments**

percent	Percent of randomly selected rows in data source
color	Colour for data points of data source

**Value**

List of plot parameters for data source

**Examples**

```
## Not run:  
gdPlotDataSourceParameters(2500)  
## End(Not run)
```

---

gdPlotParameters	<i>Specify plot parameters for generative data</i>
------------------	--

---

**Description**

Specify plot parameters for generative data passed to function `gdPlotProjection()`. When density value thresholds with assigned colors are specified generative data is drawn for density value ranges in increasing order.

**Usage**

```
gdPlotParameters(  
  percent = 10,  
  densityValueThresholds = c(),  
  densityValueColors = c("red")  
)
```

**Arguments**

<code>percent</code>	Percent of randomly selected rows in generative data
<code>densityValueThresholds</code>	Vector of density value thresholds
<code>densityValueColors</code>	Vector of colors assigned to density value thresholds. The size must be the size of <code>densityValueThresholds</code> plus one.

**Value**

List of plot parameters for generative data

**Examples**

```
## Not run:  
gdPlotParameters(50, c(0.75), c("red", "green"))  
## End(Not run)
```

---

gdPlotProjection	<i>Create an image file for generative data and data source</i>
------------------	---

---

**Description**

Create an image file containing two-dimensional projections of generative data and data source. Plot parameters for generative data and data source are passed by functions `gdPlotParameters()` and `gdPlotDataSourceParameters()`. Data points of data source are drawn above data points of generative data.

**Usage**

```
gdPlotProjection(
  imageFileName,
  title,
  columnIndices,
  generativeDataParameters = gdPlotParameters(percent = 10, densityValueThresholds = c(),
  densityValueColors = c("red")),
  dataSourceParameters = gdPlotDataSourceParameters(percent = 100, color = "blue")
)
```

**Arguments**

`imageFileName` Name of image file

`title` Title of image

`columnIndices` Vector of two column indices that are used for the two-dimensional projections. Indices refer to indices of active columns of data source.

`generativeDataParameters` Plot generative data parameters, see function `gdPlotParameters()`.

`dataSourceParameters` Plot data source parameters, see function `gdPlotDataSourceParameters()`.

**Value**

None

**Examples**

```
## Not run:
gdRead("gd.bin", "ds.bin")
gdPlotProjection("gd12ddv.png",
  "Generative Data with a Density Value Threshold for the Iris Dataset", c(1, 2),
  gdPlotParameters(250000, c(0.71), c("red", "green")),
  gdPlotDataSourceParameters(2500))
gdPlotProjection("gd34ddv.png",
  "Generative Data with a Density Value Threshold for the Iris Dataset", c(3, 4),
  gdPlotParameters(250000, c(0.71), c("red", "green")),
  gdPlotDataSourceParameters(2500))
## End(Not run)
```

---

gdRead

*Read generative data and data source*

---

**Description**

Read generative data and data source from specified files. Read in generative data and data source are accessed in `gdPlot2dProjection()`, generative data is accessed in `gdGetRow()`, `gdCalculateDensityValue()` and `gdCalculateDensityValueQuantile()`.

**Usage**

```
gdRead(generativeDataFileName, dataSourceFileName = "")
```

**Arguments**

```
generativeDataFileName  
    Name of generative data file  
dataSourceFileName  
    Name of data source file
```

**Value**

None

**Examples**

```
## Not run:  
gdRead("gd.bin", "ds.bin")  
## End(Not run)
```

---

gdWriteSubset	<i>Write subset of generative data</i>
---------------	--

---

**Description**

Write subset of randomly selected rows of generative data

**Usage**

```
gdWriteSubset(fileName, percent)
```

**Arguments**

```
fileName    Name of subset generative data file  
percent     Percent of randomly selected rows
```

**Value**

None

**Examples**

```
## Not run:  
gdRead("gd.bin")  
gdWriteSubset("gds.bin", 50)  
## End(Not run)
```

# Index

## \* package

ganGenerativeData-package, 2

dsActivateColumns, 9

dsCreateWithDataFrame, 9

dsDeactivateColumns, 10

dsGetActiveColumnNames, 10

dsGetInactiveColumnNames, 11

dsGetNumberOfRows, 11

dsGetRow, 12

dsRead, 12

dsWrite, 13

ganGenerativeData

(ganGenerativeData-package), 2

ganGenerativeData-package, 2

gdCalculateDensityValue, 14

gdCalculateDensityValueQuantile, 14

gdCalculateDensityValues, 15

gdComplete, 16

gdGenerate, 16

gdGenerateParameters, 17

gdGetNumberOfRows, 18

gdGetRow, 19

gdKNearestNeighbors, 19

gdPlotDataSourceParameters, 20

gdPlotParameters, 21

gdPlotProjection, 21

gdRead, 22

gdWriteSubset, 23