

# Package: RFQT (via r-universe)

September 26, 2024

**Type** Package

**Title** Random Forest of Q Trees

**Version** 0.1.0

**Author** Haodong Tian

**Maintainer** Haodong Tian <haodong.tian@mrc-bsu.cam.ac.uk>

**Description** Data-adaptive method for effect heterogeneity analysis or non-linear causal studies in Mendelian randomization and instrumental variables analysis.

**License** GPL (>= 2)

**Encoding** UTF-8

**LazyData** true

**Imports** stats, MendelianRandomization, tidyverse, data.table, parallel

**Suggests** dplyr

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

**RemoteUrl** <https://github.com/HDTian/RFQT>

**RemoteRef** HEAD

**RemoteSha** 1edae49a2fcba3d1f64db12c1149947d1805469a

## Contents

BootstrapTreeFitting . . . . .	2
DTfit . . . . .	5
getDat . . . . .	8
GetIndex . . . . .	9
getMSE . . . . .	11
GetNindex . . . . .	12
getPars . . . . .	13
getPredict . . . . .	15
GetTree . . . . .	16
getVI . . . . .	18
RFQTfit . . . . .	19

<b>Index</b>	<b>23</b>
--------------	-----------

---

 BootstrapTreeFitting *One Bootstrap Q Tree Fitting*


---

## Description

`BootstrapTreeFitting` fits the Q-tree for a bootstrap data of the training set (with user-specific seed) and store the cleaned results. If you do not need results but the original (dirty) Q-tree information, you may use `GetTree()`.

## Usage

```
BootstrapTreeFitting(seed=1,
                    Odat=odat,
                    Vdat=vdat,
                    honest=FALSE,
                    S=5,
                    rate=0.4,
                    SingleM=FALSE,
                    Qthreshold=3.84,
                    method='DR',
                    SoP=10,
                    howGX='SpecificGX',
                    Halve=FALSE,
                    endsize=1000,
                    const=NA)
```

## Arguments

<code>seed</code>	seed value for reproducible and traceable results
<code>Odat</code>	a data frame. This is recognized as the trainging data set.
<code>Vdat</code>	a data frame. This is recognized as the testing or validation data set.
<code>honest</code>	Logical value indicates whether to use the honest estimation style or not. Default value is <code>FALSE</code> . When honest estimation is used, the tree construction and the leaf-specific estimates will be obtained in two seperate samples of the (bootstrapped) training data.
<code>S</code>	a postive integer indicates the maximum tree depth allowed. Default value is 5 (i.e. the data set is allowed to be splitted at most 5 times). <code>S</code> is also conneted to the argument <code>endsize</code> , both of which work in a similar way. One may only restrict one of them and leave another one be a flexible value (e.g. <code>S=100</code> or <code>endsize=0</code> ).
<code>rate</code>	a value ranging from 0 to 1 indicates the proportion of the candidate covariates to be randomly cosidered in each split when construcing the Q-tree.
<code>SingleM</code>	Logical value indicates whether to use one single covariate variable at all times. Default value is <code>FALSE</code> . If <code>SingleM=TRUE</code> , the covariate used are the best variable indicated by <code>GetIndex</code> .

<code>Qthreshold</code>	A positive value used as the threshold value for the Q statistic. Default value is 3.0. One can set <code>Qthreshold=0</code> to ignore the Q-related stopping rule.
<code>method</code>	a character indicates the stratification method used for constructing a Q-tree. There are currently three methods: the Doubly-ranked method ( <code>method='DR'</code> ), the residual methods ( <code>method='Residual'</code> ) and the naive method (recognized by any character except <code>'DR'</code> and <code>'Residual'</code> ). The default method is <code>'DR'</code> .
<code>SoP</code>	a positive integer ( $\geq 2$ ) indicates the size of pre-stratum. Only applicable for the doubly-ranked stratification. Default value is 10.
<code>howGX</code>	a character indicates the way to calculate the instrument-exposure associations. Two ways are currently allowed: the instrument-exposure associations estimated separately at each stratum ( <code>howGX='SpecificGX'</code> ), or use a fixed instrument-exposure association ( <code>howGX='const'</code> )
<code>Halve</code>	Logical. Indicates whether to split the node by half:half (i.e. 5:5). Default value is <code>FALSE</code> .
<code>endsize</code>	a positive integer value indicates the minimal size of the node of Q-tree. <code>S</code> and <code>endsize</code> work in a similar way.
<code>const</code>	a value indicates the fixed instrument-exposure association value used for constructing the Q-tree. Only applicable for <code>howGX='SpecificGX'</code> .

## Details

The training data `Odat` and the testing data `Vdat` should be of the same form and be a data frame that can be recognized by the functions. The first four columns must be the individual IDs, the (one-dimensional) instrument, the exposure (should be `NA` if not available), and the outcome, respectively. The (high-dimensional) covariates are then following. Note that all the functions may be working with incorrect column ordering, but the results are misleading. Do make sure the column variable order is correct.

If simulated data is used, there should be a end column named as `true_STE`. That is, `odat>true_STE` is not `NULL`.

If the data does not contain the exposure information, the RFQT can still be fitted, where one should use `howGX='const'` with a external instrument-exposure association (usually from independent studies).

## Value

`BootstrapTreeFitting` returns a list with the following components:

<code>end_node_information</code>	the end nodes (leaves) information of the fitted Q tree, which contains the end-node index, the MR estimates, and the sample size proportion.
<code>OOB_predict</code>	the predicted effects of the OOB samples according to the Q-tree. Values will be 0 if the individual is not out of bag.
<code>v_predict</code>	the predicted effects of the testing samples according to the Q-tree.

<code>vi1</code>	the vector of the variable importance (VI) measurements for all the candidate covariates considered.
<code>vi2</code>	the vector of the variable importance (VI) measurements for all the candidate covariates considered. This VI measurements do not need the individual true effect information, therefore suitable for real application data.
<code>ts1</code>	the value of the permutation test with statistic S1.
<code>ts2</code>	the value of the permutation test with statistic S2.

When the real data is fitted, `vi1` is not applicable. When only training data is available, `v_predict` is not applicable.

## References

Burgess, S., Davies, N. M., & Thompson, S. G. (2014). "Instrumental variable analysis with a nonlinear exposure–outcome relationship". *Epidemiology (Cambridge, Mass.)*, 25(6), 877. (Residual stratification)

Tian, H., Mason, A. M., Liu, C., & Burgess, S. (2023). "Relaxing parametric assumptions for non-linear Mendelian randomization using a doubly-ranked stratification method". *PLOS Genetics*, 19(6), e1010823. (Doubly-ranked stratification)

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

#When running RFQT with multiple Q trees/bootstrap - use parallel computation
Nb<-5 #how many trees in the forest? e.g. 5 trees
cl<-makeCluster(2)# clusters/
clusterEvalQ(cl=cl , expr=library(dplyr))
clusterEvalQ(cl=cl , expr=library(MendelianRandomization) )
clusterExport( cl=cl , varlist=c( 'odat', 'vdat',
                                  'GetTree', 'GetNindex', 'GetIndex' ) )#or any other arguments
RES<-parSapply( cl , 1:Nb, BootstrapTreeFitting )
stopCluster(cl)
dim(RES)#7 Nb
##If you wish to use your own parameters rather than the default parameters, try:
#general example
user_BootstrapTreeFitting<-function(seed){
  RES<-BootstrapTreeFitting(seed,
                            honest=my.honest,
                            S=my.S,
```

```

                                JJ=my.JJ,  #or any partial of the arguments
                                rate=my.rate,
                                Qthreshold=my.Qthreshold,
                                method=my.method,
                                SoP=my.SoP,
                                howGX=my.howGX,
                                endsize=my.endsize)

    return(RES)
}
#specific exmaple
user_BootstrapTreeFitting<-function(seed){
  RES<-BootstrapTreeFitting(seed,SoP=20)
  return(RES)
}
###parallel computation -> RFQT
Nb<-5 #how many trees in the forest?
cl<-makeCluster(2)
clusterEvalQ(cl=cl , expr=library(dplyr))
clusterEvalQ(cl=cl , expr=library(MendelianRandomization) )
clusterExport( cl=cl , varlist=c( 'odat', 'vdat',
                                'GetTree', 'GetNindex', 'GetIndex' , 'BootstrapTreeFitting') )
RES<-parSapply( cl , 1:Nb, user_BootstrapTreeFitting )
stopCluster(cl)

```

---

DTfit

*Fit a Classical Single Decision Tree (DT)*


---

## Description

DTfit fits a data set (the training data) by the classical Q tree, and return the summary results. DTfit works similarly as `GetTree` and `BootstrapTreeFitting`, but not designed or used for future random forest fitting. If you wish to fit data via single tree rather than random forest, use DTfit.

## Usage

```

DTfit(Odat = odat,
      Vdat = NA,
      honest = FALSE,
      S = 5,
      rate =1,
      SingleM = FALSE,
      Qthreshold = 3.84,
      method = "DR",
      SoP = 10,
      howGX = "SpecificGX",
      Halve = FALSE,
      endsize = 1000,
      const = NA)

```

**Arguments**

<code>Odat</code>	a data frame. This is recognized as the training data set.
<code>Vdat</code>	a data frame. This is recognized as the testing or validation data set. The default value is <code>NA</code> , indicating that no testing data will be considered.
<code>honest</code>	logical value indicates whether to use the honest estimation style or not. Default value is <code>FALSE</code> . When honest estimation is used, the tree construction and the leaf-specific estimates will be obtained in two separate samples of the (bootstrapped) training data.
<code>S</code>	a positive integer indicates the maximum tree depth allowed. Default value is 5 (i.e. the data set is allowed to be splitted at most 5 times). <code>S</code> is also connected to the argument <code>endsize</code> , both of which work in a similar way. One may only restrict one of them and leave another one be a flexible value (e.g. <code>S=100</code> or <code>endsize=0</code> ).
<code>rate</code>	a value ranging from 0 to 1 indicates the proportion of the candidate covariates to be randomly considered in each split when constructing the Q-tree.
<code>SingleM</code>	logical value indicates whether to use one single covariate variable at all times. Default value is <code>FALSE</code> . If <code>SingleM=TRUE</code> , the covariate used are the best variable indicated by <code>GetIndex</code> .
<code>Qthreshold</code>	a positive value used as the threshold value for the Q statistic. Default value is 3.84. One can set <code>Qthreshold=0</code> to ignore the Q-related stopping rule.
<code>method</code>	a character indicates the stratification method used for constructing a Q-tree. There are currently three methods: the Doubly-ranked method ( <code>method='DR'</code> ), the residual methods ( <code>method='Residual'</code> ) and the naive method (recognized by any character except <code>'DR'</code> and <code>'Residual'</code> ). The default method is <code>'DR'</code> .
<code>SoP</code>	a positive integer ( $\geq 2$ ) indicates the size of pre-stratum. Only applicable for the doubly-ranked stratification. Default value is 10.
<code>howGX</code>	a character indicates the way to calculate the instrument-exposure associations. Two ways are currently allowed: the instrument-exposure associations estimated separately at each stratum ( <code>howGX='SpecificGX'</code> ), or use a fixed instrument-exposure association ( <code>howGX='const'</code> ).
<code>Halve</code>	logical. Indicates whether to split the node by half:half (i.e. 5:5). Default value is <code>FALSE</code> .
<code>endsize</code>	a positive integer value indicates the minimal size of the node of Q-tree. <code>S</code> and <code>endsize</code> work in a similar way.
<code>const</code>	a value indicates the fixed instrument-exposure association value used for constructing the Q-tree. Only applicable for <code>howGX='SpecificGX'</code> .

**Details**

Like `RFQTfit`, `DTfit` is an integrated function combining other functions, including `getPredict` and `getVI`.

The training data `Odat` and the testing data `Vdat` should be of the same form and be a data frame that can be recognized by the functions. The first four columns must be the individual IDs, the (one-dimensional) instrument, the exposure (should be NA if not available), and the outcome, respectively. The (high-dimensional) covariates are then following. Note that all the functions may be working with incorrect column ordering, but the results are misleading. DO make sure the column variable order is correct.

## Value

`DTfit` returns a list consisting of the following components:

<code>end_node_information</code>	the end nodes (leaves) information of the fitted Q tree, which contains the end-node index, the MR estimates, and the sample size proportion.
<code>v_predicted</code>	the predicted effects of the testing samples according to the Q-tree.
<code>ts1</code>	the value of the permutation test with statistic S1.
<code>ts2</code>	the value of the permutation test with statistic S2.
<code>MSE</code>	the MSE results for the testing samples (if both of the testing set and the label exist).
<code>rdat</code>	the tree fitted data.

## References

- Burgess, S., Davies, N. M., & Thompson, S. G. (2014). "Instrumental variable analysis with a nonlinear exposure–outcome relationship". *Epidemiology (Cambridge, Mass.)*, 25(6), 877. (Residual stratification)
- Tian, H., Mason, A. M., Liu, C., & Burgess, S. (2023). "Relaxing parametric assumptions for non-linear Mendelian randomization using a doubly-ranked stratification method". *PLOS Genetics*, 19(6), e1010823. (Doubly-ranked stratification)

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

DTRES<-DTfit(odat,vdat)

DTRES

#get (testing set) MSE manually
mean( ( DTRES$v_predict - vdat$true_STE )^2 )
```

getDat

*Creat a Simulated Data***Description**

getDat creates a toy data based on a certain model

**Usage**

```
getDat(N = 150000,
      Nt = 100000,
      Nc = 20,
      scenario='A',
      SoM = 0.5,
      ZXeffect = 0.5,
      Random = TRUE,
      label = TRUE,
      split = TRUE)
```

**Arguments**

<b>N</b>	a integer indicates the total sample size (training data size + testing data size)
<b>Nt</b>	a integer indicates the size of the training data.
<b>Nc</b>	a integer indicates the number of candidate variables.
<b>scenario</b>	a character indicates the model scenarios data will be simulated from. 'A' refers to scenario where no colliders. 'B' refers to the scenario that half covariates will be colliders. 'C' refers to more complicated collider scenarios.
<b>SoM</b>	a value indicates the strength of modification (i.e. how strong will the covariate modify the treatment effect). Note even if <b>SoM=0</b> , there still exists weak modification.
<b>ZXeffect</b>	a value indicates the instrument-exposure effect.
<b>Random</b>	logical. Indicates whether to use the random value for strenght modification. Default value is <b>TRUE</b> , where weak modification always exists even if <b>SoM=0</b> due to randomness.
<b>label</b>	logical. Indicates whether to added the final column of the simulated data as the true heterogenerous effect (i.e. label). Default value is <b>TRUE</b> .
<b>split</b>	logical. If <b>FALSE</b> , only one data set is returned. The default is <b>TRUE</b> , the data will be splitted into the training data and testing data.

**Details**

The data-generating model is the same as the Scenarios in the original paper, where the first 5 covariates are the true effect modifiers. The covariates in even positions (2,4,...) are the common downstream variables of the exposure and te confounders, therefore causing collider bias (when conditioning on these variables).



**Value**

`getDat()` returns a list with three kinds of toy data set

<code>whole.data</code>	the complete whole dataset
<code>training.set</code>	the training set
<code>testing.set</code>	the testing set
<code>SoM</code>	Strength of Modification
<code>modifier_vec</code>	the modification strength vector for each covariate
<code>Scenario</code>	scenario information

the training set and the testing set are of the same form that can be recognized by all the functions. The first four columns are the individual IDs, the instrument, the exposure, and the outcome, respectively. The high-dimensional (dimensions =  $N_c$ ) covariates are then following. The end column is `true_STE`, representing the individual controlled direct treatment effect (see more details in the original paper).

**Examples**

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
odat<-res$training.set #training set
vdat<-res$testing.set #testing set
```

---

 GetIndex

*Get the Best Covariate Index*


---

**Description**

`GetIndex` gives the best covariate index according to the inputted data set. The best covariate is one out of the candidate covariates considered that gives the largest Q statistic value.

**Usage**

```
GetIndex(dat_current,
         JJ,
         rate = 1,
         SpecificM = NA,
         method = "DR",
         SoP = 10,
         howGX = "SpecificGX",
         Halve = FALSE,
         const = NA)
```

**Arguments**

<code>dat_current</code>	A data frame. This data frame has the first four columns: individual IDs, the instrument, the exposure, and the outcome with the candidate covariates following.
<code>JJ</code>	A positive integer indicates the total number of the candidate covariates.
<code>rate</code>	a value ranging from 0 to 1 indicates the proportion of the candidate covariates to be randomly considered in each split when constructing the Q-tree.
<code>SpecificM</code>	a vector indicates the index of candidate variables that can be considered in splitting. The default value is <code>NA</code> , which means all candidate variables are possible to be considered.
<code>method</code>	a character indicates the stratification method used for constructing a Q-tree. There are currently three methods: the Doubly-ranked method ( <code>method='DR'</code> ), the residual methods ( <code>method='Residual'</code> ) and the naive method (recognized by any character except <code>'DR'</code> and <code>'Residual'</code> ). The default method is <code>'DR'</code> .
<code>SoP</code>	a positive integer ( $\geq 2$ ) indicates the size of pre-stratum. Only applicable for the doubly-ranked stratification. Default value is 10.
<code>howGX</code>	a character indicates the way to calculate the instrument-exposure associations. Two ways are currently allowed: the instrument-exposure associations estimated separately at each stratum ( <code>howGX='SpecificGX'</code> ), or use a fixed instrument-exposure association ( <code>howGX='const'</code> )
<code>Halve</code>	Logical. Indicates whether to split the node by half:half (i.e. 5:5). Default value is <code>FALSE</code> .
<code>const</code>	a value indicates the fixed instrument-exposure association value used for constructing the Q-tree. Only applicable for <code>howGX='SpecificGX'</code> .

**Details**

`GetIndex` helps to decide the splitting covariate at the present node when constructing a Q-tree.

**Value**

`GetIndex` returns a vector with the following three values: `'Candidate.index'`, `'Q.value'`, `'node.size'`

**`Candidate.index`**

The index (position rank) of the chosen covariate.

**`Q.value`**

The corresponding Q statistic value of the chosen covariate.

**`minimal.node.size.after.split`**

The minimal sample size of the sub-nodes after splitting.

**`split.style`**

Three values denoted by 1,2,3 representing the splitting style used for the final covariate. They are 3:7, 5:5, 7:3, respectively.

## References

Burgess, S., Davies, N. M., & Thompson, S. G. (2014). "Instrumental variable analysis with a nonlinear exposure–outcome relationship". *Epidemiology (Cambridge, Mass.)*, 25(6), 877. (Residual stratification)

Tian, H., Mason, A. M., Liu, C., & Burgess, S. (2022). "Relaxing parametric assumptions for non-linear Mendelian randomization using a doubly-ranked stratification method". *bioRxiv*, 2022-06. (Doubly-ranked stratification)

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
odat<-res$training.set #training set
vdat<-res$testing.set #testing set
GetIndex(odat,JJ=20)
```

---

getMSE *Get MSE Values for the OOB Data or the Testing Data with Fitted Q-Tree or Random Forest of Q-Trees (RFQT)*

---

## Description

getMSE calculates the MSE value for RFQT (or a Q-tree if RFQT only contains one single tree)

## Usage

```
getMSE(RES,
       indicator = 1
       )
```

## Arguments

**RES** a fitting list result from `BootstrapTreeFitting` or `RFQTfit(odat,vdat)$RES`.

**indicator** numeric indicator to judge which type of MSE should be calculated. `indicator=1` for Out-of-Bag (OOB) data and `indicator=2` for the testing data. Default value is `indicator=1`.

## Value

getMSE returns a vector of MSE with the increase of the number of Q-trees. The stable MSE values should indicate an appropriate number of Q-trees.

## Examples

```

library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
res<-getDat(label=FALSE)
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

##When running RFQT with mutiple Q trees/bootstrap - use parallel computation
Nb<-7 #how many trees in the forest? e.g. 5 trees
cl<-makeCluster(2)
clusterEvalQ(cl=cl , expr=library(dplyr))
clusterEvalQ(cl=cl , expr=library(MendelianRandomization) )
clusterExport( cl=cl , varlist=c('odat','vdat',
                                'GetTree', 'GetNindex', 'GetIndex' ) )
RES<-parSapply( cl , 1:Nb, BootstrapTreeFitting )
stopCluster(cl)
dim(RES)#7 Nb

getMSE(RES,indicator=1 )
getMSE(RES,indicator=2 )

```

---

GetNindex

*Get the Q Tree End Node Index for Chosen Samples*

---

## Description

GetNindex helps to match the Q-tree end node for any inputed samples, and therefore enables to make further analysis like effect prediction for these inputed samples.

## Usage

```

GetNindex(M,
          rdat,
          S = NA)

```

## Arguments

M	A matrix or data frame contains the candidate covariate information.
rdat	A data frame. It is the Q-tree informaton result and obtained by <code>rdat&lt;-GetTree()</code> .
S	A positive integer indicates the max depth of each sample will explore along the tree. This number should be consistent with the max depth of the fitted Q tree (also relected in <code>rdat</code> ). The default valus if NA, which means the code will automatically use the correct S value.

## Details

The inputted data set `M` must contain the same candidate covariate information (with the same order) as the training data for the Q-tree. In addition, the first four columns (individual IDs, the instrument, the exposure, the outcome) need to be removed, so that the first column starts from the first candidate covariate.

## Value

`GetNindex` returns a vector with the same length of the row number of the inputted data `M`. Each element represents the end node (i.e. leaf) index according to the reference Q-tree (the tree information was stored within `rdat`).

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

rdat<-GetTree(odat)

vdat_Nindex<-GetNindex(vdat[,5:24] ,rdat )
#the M colnumber and order should be same as the training data

#may use another independent data (estimation data) to
#calculate the endnode-specific IV/MR estimates

#for example:
#let odat halves into trdata (tree data) and estdata (estimation data);
#also we can have rdat<-GetTree(trdata)

trdat<-odat[1: (nrow(odat)/2) ,]#tree data
estdat<-odat[(nrow(odat)/2+1):nrow(odat) ,] #estimation data

rdat<-GetTree(trdat)#fitted Q tree

estdat$Nindex<-GetNindex(estdat[,5:24] ,rdat )
#rdat contains the tree information (i.e. decision rule)
```

## Description

`getPars` returns the information of the (hyper-)parameters relevant to the Q-tree and RFQT fitting.

## Usage

```
getPars(empty.argument)
```

## Arguments

`empty.argument`  
No arguments are needed.

## Value

`getPars()` returns the following information: the total number of candidate covariates, the training data size, the testing data size, the stratification method, the size of pre-stratum (if applicable), the proportion of the covariates randomly considered in each node split, the maximal tree depth allowed, the way to calculate the instrument-exposure associations, the fixed instrument-exposure association level (if applicable), the minimal node size, the threshold value of the Q statistic.

If the variable is not defined, it will return 'Not defined' with a value in brackets representing the default value for this variable when fitting RFQT.

Note that only the variable defined in the global environment will be considered. Even if one defines these variables, the functions like `BootstrapTreeFitting` and `RFQTfit` will use the default value if (s)he does not modify the function arguments.

See the original paper for the details of each parameter information.

## Author(s)

Haodong Tian

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

res<-getDat()
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

method<-'DR'; SoP<-20; rate<-2/5

getPars()
```

---

getPredict	<i>Get predict effects for the OOB Data or the Testing Data with Fitted Random Forest of Q-Trees (RFQT)</i>
------------	---

---

## Description

getPredict calculates the predicted individual treatment effect estimates for RFQT, either for the OOB data or the testing data.

## Usage

```
getPredict(RES,
           indicator = 1
           )
```

## Arguments

RES	a fitting list result from <code>BootstrapTreeFitting</code> or <code>RFQTfit(odat, vdat)\$RES</code> .
indicator	numeric indicator to judge which type of predicted effect should be presented. <code>indicator=1</code> for Out-of-Bag (OOB) data and <code>indicator=2</code> for th testing data. Defalut value is <code>indicator=1</code> .

## Details

See the relevant part of the original paper for more details of how to predict an individual effect according to a fitted Q tree (therefore a RFQT).

## Value

getPredict returns a matrix where the row corresponds to the sample size (either the OOB data or the testing data, depending on the indicator used) and the column coresponds to the number of Q-trees. Each column represents the individual predicted effects with the present number of Q trees.

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the deflault setting: scenario='A' and SoM=0.5
res<-getDat(label=FALSE)
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

##When running RFQT with mutiple Q trees/bootstrap - use parallel computation
Nb<-7 #how many trees in the forest? e.g. 5 trees
```

```

cl<-makeCluster(2)
clusterEvalQ(cl=cl , expr=library(dplyr))
clusterEvalQ(cl=cl , expr=library(MendelianRandomization) )
clusterExport( cl=cl , varlist=c('odat','vdat',
                                'GetTree', 'GetNindex', 'GetIndex' ) )
RES<-parSapply( cl , 1:Nb, BootstrapTreeFitting )
stopCluster(cl)
dim(RES)#7 Nb

predict_matrix<-getPredict(RES,indicator=1 )
predict_matrix<-getPredict(RES,indicator=2 )

dim(predict_matrix)
###View(predict_matrix)

```

---

GetTree

*Fit and Get a Q Tree*


---

## Description

**GetTree** fits a data and construct a Q-tree. The fitted data is directly used to build tree (so no bootstrap). The tree information will be stored but not cleaned for further analysis. You may need **BootstrapTreeFitting** rather than this function.

## Usage

```

GetTree(dat,
        S = 5,
        Qthreshold = 3.84,
        rate = 1,
        SpecificM = NA,
        method = "DR",
        SoP = 10,
        howGX = "SpecificGX",
        Halve = FALSE,
        const = NA,
        endsize = 1000)

```

## Arguments

<b>dat</b>	a data frame. This data (without any bootstrap) will be fitted to construct a Q-tree.
<b>S</b>	a positive integer indicates the maximum tree depth allowed. Default value is 5 (i.e. the data set is allowed to be splitted at most 5 times). <b>S</b> is also conneted to the argument <b>endsize</b> , both of which work in a similar way. One may only restrict one of them and leave another one be a flexible value (e.g. <b>S=100</b> or <b>endsize=0</b> ).



<code>Qthreshold</code>	A positive value used as the threshold value for the Q statistic. Default value is 3.0.
<code>rate</code>	a value ranging from 0 to 1 indicates the proportion of the candidate covariates to be randomly considered in each split when constructing the Q-tree. When <code>rate=1</code> , it means all the candidate covariates will be considered.
<code>SpecificM</code>	a vector indicates the index of candidate variables that can be considered in splitting. The default value is <code>NA</code> , which means all candidate variables are possible to be considered.
<code>method</code>	a character indicates the stratification method used for constructing a Q-tree. There are currently three methods: the Doubly-ranked method ( <code>method='DR'</code> ), the residual methods ( <code>method='Residual'</code> ) and the naive method (recognized by any character except <code>'DR'</code> and <code>'Residual'</code> ). The default method is <code>'DR'</code> .
<code>SoP</code>	a positive integer ( $\geq 2$ ) indicates the size of pre-stratum. Only applicable for the doubly-ranked stratification. Default value is 10.
<code>howGX</code>	a character indicates the way to calculate the instrument-exposure associations. Two ways are currently allowed: the instrument-exposure associations estimated separately at each stratum ( <code>howGX='SpecificGX'</code> ), or use a fixed instrument-exposure association ( <code>howGX='const'</code> )
<code>Halve</code>	Logical. Indicates whether to split the node by half:half (i.e. 5:5). Default value is <code>FALSE</code> .
<code>const</code>	a value indicates the fixed instrument-exposure association value used for constructing the Q-tree. Only applicable for <code>howGX='SpecificGX'</code> .
<code>endsize</code>	a positive integer value indicates the minimal size of the node of Q-tree allowed. One can also control the endsize via <code>S</code> .

## Details

`dat` must be a data frame that can be recognized by the functions. The first four columns must be the individual IDs, the (one-dimensional) instrument, the exposure (should be `NA` if not available), and the outcome, respectively. The (high-dimensional) covariates are then following. Note that all the functions may be working with incorrect column ordering, but the results are misleading. DO make sure the column variable order is correct.

If simulated data is used, there should be a end column named as `true_STE`. That is, `dat>true_STE` is not `NULL`.

## Value

`GetTree` returns a data set with the same row dimension as the inputted data `dat`. The returned data is 'augmented' with tree information including the tree end node index, and the splitting information in each split.

## References

Burgess, S., Davies, N. M., & Thompson, S. G. (2014). "Instrumental variable analysis with a nonlinear exposure-outcome relationship". *Epidemiology (Cambridge, Mass.)*, 25(6), 877. (Residual stratification)

Tian, H., Mason, A. M., Liu, C., & Burgess, S. (2023). "Relaxing parametric assumptions for non-linear Mendelian randomization using a doubly-ranked stratification method". *PLOS Genetics*, 19(6), e1010823. (Doubly-ranked stratification)

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

res<-getDat()
odat<-res$training.set #training set

tree_result<-GetTree(odat)
```

---

getVI	<i>Get Variable Importance Measurements with Fitted Q-Tree or Random Forest of Q-Trees (RFQT)</i>
-------	---

---

## Description

getVI calculates the Variable Importance (VI) measurements for a random forest of Q-trees (RFQT).

## Usage

```
getVI(RES,
      VItype = 2
      )
```

## Arguments

RES	a fitting list result from <code>BootstrapTreeFitting</code> or <code>RFQTfit(odat,vdat)\$RES</code> .
VItype	numeric indicator to judge which type of VI measurement should be calculated. <code>VItype=1</code> for the classical way with the known individual treatment effects (i.e. label) and <code>VItype=2</code> for case that the individual treatment effects are unavailable (as in most real application). Default value is <code>VItype=2</code> .

## Details

See the VI measurement part of the original paper for more detailed VI introduction and the algorithm.

## Value

getVI returns a vector indicating the variable importance (VI) measurements for all the candidate covariates (with the same ordering of the candidate covariates to the original fitting data).

## Examples

```

library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
res<-getDat(label=FALSE)
odat<-res$training.set #training set
vdat<-res$testing.set #testing set

##When running RFQT with mutiple Q trees/bootstrap - use parallel computation
Nb<-7 #how many trees in the forest? e.g. 5 trees
cl<-makeCluster(2)
clusterEvalQ(cl=cl , expr=library(dplyr))
clusterEvalQ(cl=cl , expr=library(MendelianRandomization) )
clusterExport( cl=cl , varlist=c('odat','vdat',
                                'GetTree', 'GetNindex', 'GetIndex' ) )
RES<-parSapply( cl , 1:Nb, BootstrapTreeFitting )
stopCluster(cl)
dim(RES)#7 Nb

getVI(RES,Vitype=1 )
getVI(RES,Vitype=2 )

```

---

RFQFit

*Fit a Random Forest of Q Trees (RFQT)*


---

## Description

RFQFit fits a data set (the training data) by the random forest of Q trees (RFQT), and return the summary results.

## Usage

```

RFQFit(odat,
       vdat = NA,
       Nb = 5,
       S = 5,
       honest = FALSE,
       rate = 0.4,
       SingleM = FALSE,
       Qthreshold = 3.84,
       method = "DR",
       SoP = 10,
       howGX = "SpecificGX",
       Halve = FALSE,

```

```

endsize = 1000,
const = NA,
Cores = NA,
trackfile=NA
)

```

### Arguments

<code>odat</code>	a data frame. This is recognized as the training data set.
<code>vdat</code>	a data frame. This is recognized as the testing or validation data set. The default value is <code>NA</code> , indicating that no testing data will be considered.
<code>Nb</code>	a positive integer indicates the number of bootstrap (i.e. the number of Q-trees)
<code>S</code>	a positive integer indicates the maximum tree depth allowed. Default value is 5 (i.e. the data set is allowed to be splitted at most 5 times). <code>S</code> is also connected to the argument <code>endsize</code> , both of which work in a similar way. One may only restrict one of them and leave another one be a flexible value (e.g. <code>S=100</code> or <code>endsize=0</code> ).
<code>honest</code>	logical value indicates whether to use the honest estimation style or not. Default value is <code>FALSE</code> . When honest estimation is used, the tree construction and the leaf-specific estimates will be obtained in two separate samples of the (bootstrapped) training data.
<code>rate</code>	a value ranging from 0 to 1 indicates the proportion of the candidate covariates to be randomly considered in each split when constructing the Q-tree.
<code>SingleM</code>	logical value indicates whether to use one single covariate variable at all times. Default value is <code>FALSE</code> . If <code>SingleM=TRUE</code> , the covariate used are the best variable indicated by <code>GetIndex</code> .
<code>Qthreshold</code>	a positive value used as the threshold value for the Q statistic. Default value is 3.84. One can set <code>Qthreshold=0</code> to ignore the Q-related stopping rule.
<code>method</code>	a character indicates the stratification method used for constructing a Q-tree. There are currently three methods: the Doubly-ranked method ( <code>method='DR'</code> ), the residual methods ( <code>method='Residual'</code> ) and the naive method (recognized by any character except <code>'DR'</code> and <code>'Residual'</code> ). The default method is <code>'DR'</code> .
<code>SoP</code>	a positive integer ( $\geq 2$ ) indicates the size of pre-stratum. Only applicable for the doubly-ranked stratification. Default value is 10.
<code>howGX</code>	a character indicates the way to calculate the instrument-exposure associations. Two ways are currently allowed: the instrument-exposure associations estimated separately at each stratum ( <code>howGX='SpecificGX'</code> ), or use a fixed instrument-exposure association ( <code>howGX='const'</code> ).
<code>Halve</code>	logical. Indicates whether to split the node by half:half (i.e. 5:5). Default value is <code>FALSE</code> .
<code>endsize</code>	a positive integer value indicates the minimal size of the node of Q-tree. <code>S</code> and <code>endsize</code> work in a similar way.

<code>const</code>	a value indicates the fixed instrument-exposure association value used for constructing the Q-tree. Only applicable for <code>howGX='SpecificGX'</code> .
<code>Cores</code>	a positive integer indicates the number of cores should be used when running random forest in parallel. Default value is the number of CPU cores minus one. You can check the number of CPU cores via <code>detectCores()</code> .
<code>trackfile</code>	the path string indicates the file to store the running information in parallel.

## Details

RFQTfit is an integrated function combining other functions, including `GetTree`, `BootstrapTreeFitting`, `getPredict`, `getVI`, etc. The fitting automatically uses parallel computation.

The training data `Odat` and the testing data `Vdat` should be of the same form and be a data frame that can be recognized by the functions. The first four columns must be the individual IDs, the (one-dimensional) instrument, the exposure (should be `NA` if not available), and the outcome, respectively. The (high-dimensional) covariates are then following. Note that all the functions may be working with incorrect column ordering, but the results are misleading. DO make sure the column variable order is correct.

## Value

RFQTfit shows all the parameter information used for random forest fitting in the beginning, and returns in the end a list consisting of the following components:

<code>RES</code>	The direct fitting results of each Q-tree, which contains <code>end_node_information</code> , <code>OOB_predict</code> , <code>v_predict</code> , <code>vi1</code> , <code>vi2</code> , <code>ts1</code> and <code>ts2</code> (same as the results of <code>BootstrapTreeFitting</code> ) for each Q-tree. Note that <code>RES</code> is multi-dimensional. If one wishes to obtain the individual predicted effects of the testing set according to <code>i</code> -th Q-tree, should run <code>RES[3,i]\$v_predict</code> .
<code>MSE_OOB</code>	the <code>Nb</code> -length vector of MSE results for OOB samples. If one determine the number of Q-trees to be <code>Nb</code> , simply use the final vlue of <code>MSE_OOB</code> as the MSE.
<code>MSE_test</code>	the <code>Nb</code> -length vector of MSE results for test samples (if the testing set exists). If one determine the number of Q-trees to be <code>Nb</code> , simply use the final vlue of <code>MSE_test</code> as the MSE.
<code>Predict_OOB</code>	the matrix of predicted effects for OOB samples, where the row corresponds to the sample size (the testing data) and the column coresponds to the number of Q-trees. Each column represents the individual predicted effects with the present number of Q trees. If one determine the number of Q-trees to be <code>Nb</code> , simply use the final column of <code>Predict_OOB</code> as the individual predicted effect vector.
<code>Predict_test</code>	the matrix of predicted effects for testing set, where the row corresponds to the sample size (the testing data) and the column coresponds to the number of Q-trees. Each column represents the individual predicted effects with the present number of Q trees. If one determine the number of Q-trees to be <code>Nb</code> , simply use the final column of <code>Predict_test</code> as the individual predicted effect vector.

VI1	the vector indicates the type I (i.e. with labels) variable importance (VI) measurements for all the candidate covariates (with the same ordering of the candidate covariates to the original fitting data).
VI2	the vector indicates the type II (i.e. without labels) variable importance (VI) measurements for all the candidate covariates (with the same ordering of the candidate covariates to the original fitting data).
TS	the vector of the permutation test statistic values. The two values represent the values of TS1 and TS2, respectively.

## References

Burgess, S., Davies, N. M., & Thompson, S. G. (2014). "Instrumental variable analysis with a nonlinear exposure–outcome relationship". *Epidemiology (Cambridge, Mass.)*, 25(6), 877. (Residual stratification)

Tian, H., Mason, A. M., Liu, C., & Burgess, S. (2023). "Relaxing parametric assumptions for non-linear Mendelian randomization using a doubly-ranked stratification method". *PLOS Genetics*, 19(6), e1010823. (Doubly-ranked stratification)

## Examples

```
library(MendelianRandomization)
library(tidyverse)
library(data.table)
library(parallel)

set.seed(60)
res<-getDat() #simulated data #the default setting: scenario='A' and SoM=0.5
odat<-res$training.set
vdat<-res$testing.set

###try: ALLRES<-RFQTfit(odat)

###try: ALLRES<-RFQTfit(odat,vdat,Nb=200)
```

# Index

BootstrapTreeFitting, [2](#)

DTfit, [5](#)

getDat, [8](#)

GetIndex, [9](#)

getMSE, [11](#)

GetNindex, [12](#)

getPars, [13](#)

getPredict, [15](#)

GetTree, [16](#)

getVI, [18](#)

RFQTfit, [19](#)